

МИНИСТЕРСТВО ОБРАЗОВАНИЯ СТАВРОПОЛЬСКОГО КРАЯ
КРАЕВОЙ ЦЕНТР РАЗВИТИЯ ТВОРЧЕСТВА ДЕТЕЙ И ЮНОШЕСТВА
имени Ю.А. Гагарина

МАЛАЯ ТЕХНИЧЕСКАЯ АКАДЕМИЯ

Лекционный материал по курсу

"ИНФОРМАЦИОННЫЕ ТЕХНОЛОГИИ"

1 год обучения (1 семестр)



г. Ставрополь, 2018 г.

Составитель: Литвинова Наталья Николаевна – кандидат педагогических наук.

Данные материалы помогут освоить базовые понятия и термины информатики, а также приобрести навыки работы на ПК.

Адрес: 355008, г. Ставрополь, ул. Комсомольская, 65, кабинет 173.

Телефон: (8652) 26-83-88.

E-mail: mta_stav@mail.ru

Краевой Центр развития творчества детей и юношества имени Ю.А.Гагарина
г. Ставрополь, 2018 год





1.1. Что такое информатика

Термин "**информатика**" (франц. *informatique*) происходит от французских слов *information* (информация) и *automatique* (автоматика) и дословно означает "информационная автоматика". Широко распространён также англоязычный вариант этого термина — "*Computer science*", что означает буквально "компьютерная наука".

Информатика — это основанная на использовании компьютерной техники дисциплина, изучающая структуру и общие свойства информации, а также закономерности и методы её создания, хранения, поиска, преобразования, передачи и применения в различных сферах человеческой деятельности.

В 1978 году Международный научный конгресс официально закрепил за понятием "информатика" области, связанные с разработкой, созданием, использованием и материально-техническим обслуживанием систем обработки информации, включая компьютеры и их программное обеспечение, а также организационные, коммерческие, административные и социально-политические аспекты компьютеризации — массового внедрения компьютерной техники во все области жизни людей.

Таким образом, информатика базируется на компьютерной технике и немыслима без нее.

Информатика — комплексная научная дисциплина с широчайшим диапазоном применения. Её приоритетные направления:

- разработка вычислительных систем и программного обеспечения;
- теория информации, изучающая процессы, связанные с передачей, приёмом, преобразованием и хранением информации;
- математическое моделирование, методы вычислительной и прикладной математики и их применение к фундаментальным и прикладным исследованиям в различных областях знаний;
- методы искусственного интеллекта, моделирующие методы логического и аналитического мышления в интеллектуальной деятельности человека (логический вывод, обучение, понимание речи, визуальное восприятие, игры и др.);
- системный анализ, изучающий методологические средства, используемые для подготовки и обоснования решений по сложным проблемам различного характера;
- биоинформатика, изучающая информационные процессы в биологических системах;
- социальная информатика, изучающая процессы информатизации общества;
- методы машинной графики, анимации, средства мультимедиа;
- телекоммуникационные системы и сети, в том числе, глобальные компьютерные сети, объединяющие всё человечество в единое информационное сообщество;
- разнообразные приложения, охватывающие производство, науку, образование, медицину, торговлю, сельское хозяйство и все другие виды хозяйственной и общественной деятельности.

В информатике выделяет три неразрывно и существенно связанные части — технические средства, программные и алгоритмические.

Технические средства, или аппаратура компьютеров, в английском языке обозначаются словом **Hardware**, которое буквально переводится как "твёрдые изделия".

Для обозначения программных средств, под которыми понимается совокупность всех программ, используемых компьютерами, и область деятельности по их созданию и применению, используется слово **Software** (буквально — "мягкие изделия"), которое подчеркивает равнозначность самой машины и программного обеспечения, а также способность программного обеспечения модифицироваться, приспосабливаться и развиваться.

Программированию задачи всегда предшествует разработка способа ее решения в виде последовательности действий, ведущих от исходных данных к искомому результату, иными словами, разработка алгоритма решения задачи. Для обозначения части информатики, связанной с разработкой алгоритмов и изучением методов и приемов их построения, применяют термин **Brainware** (англ. *brain* — интеллект).

Роль информатики в развитии общества чрезвычайно велика. С ней связано начало революции в области накопления, передачи и обработки информации. Эта революция, следующая

за революциями в овладении веществом и энергией, затрагивает и коренным образом преобразует не только сферу материального производства, но и интеллектуальную, духовную сферы жизни.

Прогрессивное увеличение возможностей компьютерной техники, развитие информационных сетей, создание новых информационных технологий приводят к значительным изменениям во всех сферах общества: в производстве, науке, образовании, медицине и т.д.



1.2. Что такое информация

Термин "*информация*" происходит от латинского слова "*informatio*", что означает сведения, разъяснения, изложение. В настоящее время понятие информация, получает различные смысловые наполнения в различных отраслях человеческой деятельности:

- в обиходе информацией называют любые данные или сведения, которые кого-либо интересуют. "Информировать" в этом смысле означает "сообщить нечто, неизвестное раньше";

- в технике под информацией понимают сообщения, передаваемые в форме знаков или сигналов;

- в кибернетике под информацией понимает ту часть знаний, которая используется для ориентирования, активного действия, управления, т.е. в целях сохранения, совершенствования, развития системы.

Клод Шеннон, американский учёный, заложивший основы теории информации — науки, изучающей процессы, связанные с передачей, приёмом, преобразованием и хранением информации, — рассматривает информацию как снятую неопределенность наших знаний о чем-то.

Современное научное представление об информации очень точно сформулировал Норберт Винер, "отец" кибернетики. Информация — это обозначение содержания, полученного из внешнего мира в процессе нашего приспособления к нему и приспособления к нему наших чувств.

В случаях, когда говорят об автоматизированной работе с информацией посредством каких-либо технических устройств, обычно в первую очередь интересуются не содержанием сообщения, а тем, сколько символов это сообщение содержит.

Применительно к компьютерной обработке данных под информацией понимают некоторую последовательность символических обозначений (букв, цифр, закодированных графических образов и звуков и т.п.), несущую смысловую нагрузку и представленную в понятном компьютеру виде. Каждый новый символ в такой последовательности символов увеличивает информационный объём сообщения.



1.3. В каком виде существует информация

Информация может существовать в виде:

- текстов, рисунков, чертежей, фотографий;

- световых или звуковых сигналов;

- радиоволн;

- электрических и нервных импульсов;

- магнитных записей;

- жестов и мимики;

- запахов и вкусовых ощущений;

- хромосом, посредством которых передаются по наследству признаки и свойства организмов и т.д.

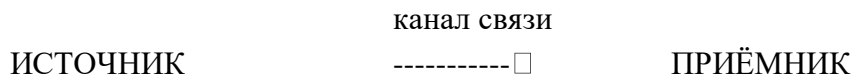
Предметы, процессы, явления материального или нематериального свойства, рассматриваемые с точки зрения их информационных свойств, называются информационными объектами.



1.4. Как передаётся информация

Информация передаётся в форме сообщений от некоторого источника информации к её приёмнику посредством канала связи между ними. Источник посылает передаваемое сообщение,

которое кодируется в передаваемый сигнал. Этот сигнал посылается по каналу связи. В результате в приёмнике появляется принимаемый сигнал, который декодируется и становится принимаемым сообщением.



Примеры:

Сообщение, содержащее информацию о прогнозе погоды, передаётся приёмнику (телезрителю) от источника — специалиста-метеоролога посредством канала связи — телевизионной передающей аппаратуры и телевизора.

Живое существо своими органами чувств (глаз, ухо, кожа, язык и т.д.) воспринимает информацию из внешнего мира, перерабатывает её в определенную последовательность нервных импульсов, передает импульсы по нервным волокнам, хранит в памяти в виде состояния нейронных структур мозга, воспроизводит в виде звуковых сигналов, движений и т.п., использует в процессе своей жизнедеятельности.

Передача информации по каналам связи часто сопровождается воздействием помех, вызывающих искажение и потерю информации.



1.5. Как измеряется количество информации

В настоящее время получили распространение подходы к определению понятия "количество информации", основанные на том, что информацию, содержащуюся в сообщении, можно нестрого трактовать в смысле её новизны или, иначе, уменьшения неопределённости наших знаний об объекте. Эти подходы используют математические понятия вероятности и логарифма.

*Подходы к определению количества информации.
Формулы Хартли и Шеннона.*

Американский инженер Р. Хартли в 1928 г. процесс получения информации рассматривал как выбор одного сообщения из конечного наперёд заданного множества из N равновероятных сообщений, а количество информации I , содержащееся в выбранном сообщении, определял как двоичный логарифм N .

Формула Хартли: $I = \log_2 N$

Допустим, нужно угадать одно число из набора чисел от единицы до ста. По формуле Хартли можно вычислить, какое количество информации для этого требуется: $I = \log_2 100 = 6,644$. Таким образом, сообщение о верно угаданном числе содержит количество информации, приблизительно равное 6,644 единицы информации.

Клод Шеннон предложил в 1948 г. другую формулу определения количества информации, учитывающую возможную неодинаковую вероятность сообщений в наборе.

Формула Шеннона: $I = - (p_1 \log_2 p_1 + p_2 \log_2 p_2 + \dots + p_N \log_2 p_N)$,
где p_i — вероятность того, что именно i -е сообщение выделено
в наборе из N сообщений.

Легко заметить, что если вероятности p_1, \dots, p_N равны, то каждая из них равна $1/N$, и формула Шеннона превращается в формулу Хартли.

Помимо двух рассмотренных подходов к определению количества информации, существуют и другие. Важно помнить, что любые теоретические результаты применимы лишь к определённым кругу случаев, очерченному первоначальными допущениями.

В качестве единицы информации Клод Шеннон предложил принять один **бит** (англ. *bit* — *binary digit* — двоичная цифра).

Бит в теории информации — количество информации, необходимое для различения двух равновероятных сообщений (типа "орел"- "решка", "чет"- "нечет" и т.п.). В вычислительной технике битом называют наименьшую "порцию" памяти компьютера, необходимую для хранения одного из двух знаков "0" и "1", используемых для внутримашинного представления данных и команд.

Бит — слишком мелкая единица измерения. На практике чаще применяется более крупная единица — байт, равная восьми битам. Именно восемь битов требуется для того, чтобы закодировать любой из 256 символов алфавита клавиатуры компьютера ($256=2^8$).

Широко используются также ещё более крупные производные единицы информации:

1 Килобайт (Кбайт) = 1024 байт = 2^{10} байт,

1 Мегабайт (Мбайт) = 1024 Кбайт = 2^{20} байт,

1 Гигабайт (Гбайт) = 1024 Мбайт = 2^{30} байт.

В последнее время в связи с увеличением объёмов обрабатываемой информации входят в употребление такие производные единицы, как:

1 Терабайт (Тбайт) = 1024 Гбайт = 2^{40} байт,

1 Петабайт (Пбайт) = 1024 Тбайт = 2^{50} байт.

За единицу информации можно было бы выбрать количество информации, необходимое для различения, например, десяти равновероятных сообщений. Это будет не двоичная (бит), а десятичная (дит) единица информации.



1.6. Что можно делать с информацией

Информацию можно:

создавать;	формализовать;	собирать;
передавать;	распространять;	хранить;
воспринимать;	преобразовывать;	искать;
использовать;	комбинировать;	измерять;
запоминать;	обрабатывать;	разрушать;
принимать;	делить на части;	и др.
копировать;	упрощать;	

Все эти процессы, связанные с определенными операциями над информацией, называются *информационными процессами*.



1.7. Какими свойствами обладает информация

Свойства информации:

достоверность;	понятность;
полнота;	доступность;
ценность;	краткость;
своевременность;	и др.

Информация достоверна, если она отражает истинное положение дел. Недостоверная информация может привести к неправильному пониманию или принятию неправильных решений.

Достоверная информация со временем может стать недостоверной, так как она обладает свойством устаревать, то есть перестаёт отражать истинное положение дел.

Информация полна, если её достаточно для понимания и принятия решений. Как неполная, так и избыточная информация сдерживает принятие решений или может повлечь ошибки.

Точность информации определяется степенью ее близости к реальному состоянию объекта, процесса, явления и т.п.

Ценность информации зависит от того, насколько она важна для решения задачи, а также от того, насколько в дальнейшем она найдёт применение в каких-либо видах деятельности человека. Только своевременно полученная информация может принести ожидаемую пользу. Одинаково нежелательны как преждевременная подача информации (когда она ещё не может быть усвоена), так и её задержка.

Если ценная и своевременная информация выражена непонятным образом, она может стать бесполезной.

Информация становится понятной, если она выражена языком, на котором говорят те, кому предназначена эта информация.

Информация должна преподноситься в доступной (по уровню восприятия) форме. Поэтому одни и те же вопросы по-разному излагаются в школьных учебниках и научных изданиях.

Информацию по одному и тому же вопросу можно изложить кратко (сжато, без несущественных деталей) или пространно (подробно, многословно). Краткость информации необходима в справочниках, энциклопедиях, учебниках, всевозможных инструкциях.



1.8. Что такое обработка информации

Обработка информации — получение одних информационных объектов из других информационных объектов путем выполнения некоторых алгоритмов

Обработка является одной из основных операций, выполняемых над информацией, и главным средством увеличения объёма и разнообразия информации.

Средства обработки информации — это всевозможные устройства и системы, созданные человеком, и в первую очередь, компьютер — универсальная машина для обработки информации.

Компьютеры обрабатывают информацию путем выполнения некоторых алгоритмов.



1.9. Что такое информационные ресурсы и информационные технологии

Информационные ресурсы — это идеи человечества и указания по их реализации, накопленные в форме, позволяющей их воспроизводство.

Это книги, статьи, патенты, диссертации, научно-исследовательская и опытно-конструкторская документация, технические переводы, данные о передовом производственном опыте и др.

Информационные ресурсы (в отличие от всех других видов ресурсов — трудовых, энергетических, минеральных и т.д.) тем быстрее растут, чем больше их расходуют.

Информационная технология — это совокупность методов и устройств, используемых людьми для обработки информации.

Человечество занималось обработкой информации тысячи лет. Первые информационные технологии основывались на использовании счётов и письменности. Около шестидесяти лет назад началось исключительно быстрое развитие этих технологий, что в первую очередь связано с появлением компьютеров.

В настоящее время термин "информационная технология" употребляется в связи с использованием компьютеров для обработки информации. Информационные технологии охватывают всю вычислительную технику и технику связи и, отчасти, — бытовую электронику, телевидение и радиовещание.



1.10. Что понимают под информатизацией общества

Информатизация общества — организованный социально-экономический и научно-технический процесс создания оптимальных условий для удовлетворения информационных потребностей и реализации прав граждан, органов государственной власти, органов местного самоуправления организаций, общественных объединений на основе формирования и использования информационных ресурсов.

Цель информатизации — улучшение качества жизни людей за счет увеличения производительности и облегчения условий их труда.

Информатизация — это сложный социальный процесс, связанный со значительными изменениями в образе жизни населения.



Общие принципы организации и работы компьютеров

2.1. Что такое компьютер

Компьютер (англ. computer — вычислитель) представляет собой программируемое электронное устройство, способное обрабатывать данные и производить вычисления, а также выполнять другие задачи манипулирования символами.

Существует два основных класса компьютеров:

- цифровые компьютеры, обрабатывающие данные в виде двоичных кодов;
- аналоговые компьютеры, обрабатывающие непрерывно меняющиеся физические величины (электрическое напряжение, время и т.д.), которые являются аналогами вычисляемых величин.

В настоящее время подавляющее большинство компьютеров являются цифровыми.

Основу компьютеров образует аппаратура (HardWare), построенная, в основном, с использованием электронных и электромеханических элементов и устройств. Принцип действия компьютеров состоит в выполнении программ (SoftWare) — заранее заданных, четко определённых последовательностей арифметических, логических и других операций.

Любая компьютерная программа представляет собой последовательность отдельных команд.

Команда — это описание операции, которую должен выполнить компьютер. Как правило, у команды есть свой код (условное обозначение), исходные данные (операнды) и результат.

Например, у команды "сложить два числа" операндами являются слагаемые, а результатом — их сумма. А у команды "стоп" операндов нет, а результатом является прекращение работы программы.

Совокупность команд, выполняемых данным компьютером, называется системой команд этого компьютера.



2.2. Как устроен компьютер

Разнообразие современных компьютеров очень велико. Но их структуры основаны на общих логических принципах, позволяющих выделить в любом компьютере следующие главные устройства:

- ✚ Память (запоминающее устройство, ЗУ), состоящую из перенумерованных ячеек;
- ✚ Процессор, включающий в себя устройство управления (УУ) и арифметико-логическое устройство (АЛУ);
- ✚ Устройство ввода;
- ✚ Устройство вывода.

Эти устройства соединены каналами связи, по которым передается информация.

Основные устройства компьютера и связи между ними представлены на схеме (рис. 2.1). Жирными стрелками показаны пути и направления движения информации, а простыми стрелками — пути и направления передачи управляющих сигналов.

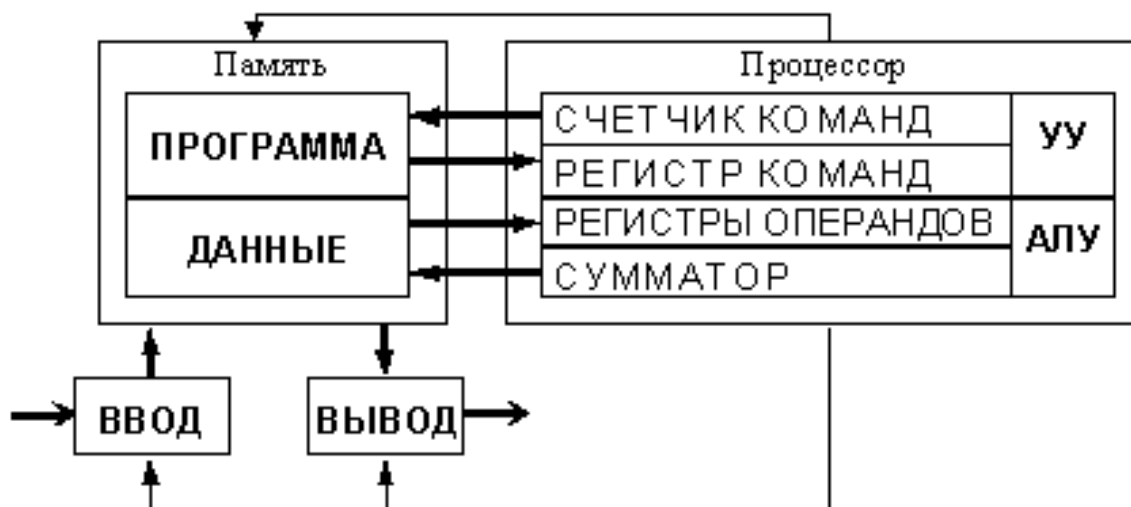


Рис. 2.1. Общая схема компьютера

Функции памяти:

- приём информации из других устройств;
- запоминание информации;
- выдача информации по запросу в другие устройства машины.

Функции процессора:

- обработка данных по заданной программе путем выполнения арифметических и логических операций;
- программное управление работой устройств компьютера.

Та часть процессора, которая выполняет команды, называется арифметико-логическим устройством (АЛУ), а другая его часть, выполняющая функции управления устройствами, называется устройством управления (УУ).

Обычно эти два устройства выделяются чисто условно, конструктивно они не разделены.

В составе процессора имеется ряд специализированных дополнительных ячеек памяти, называемых регистрами.

Регистр выполняет функцию кратковременного хранения числа или команды. Над содержимым некоторых регистров специальные электронные схемы могут выполнять некоторые манипуляции. Например, "вырезать" отдельные части команды для последующего их использования или выполнять определенные арифметические операции над числами.

Основным элементом регистра является электронная схема, называемая триггером, которая способна хранить одну двоичную цифру (разряд двоичного кода).

Регистр представляет собой совокупность триггеров, связанных друг с другом определенным образом общей системой управления.

Существует несколько типов регистров, отличающихся видом выполняемых операций. Некоторые важные регистры имеют свои названия, например:

- ✚ сумматор— регистр АЛУ, участвующий в выполнении каждой операции;
- ✚ счетчик команд— регистр УУ, содержимое которого соответствует адресу очередной выполняемой команды; служит для автоматической выборки программы из последовательных ячеек памяти;
- ✚ регистр команд — регистр УУ для хранения кода команды на период времени, необходимый для ее выполнения. Часть его разрядов используется для хранения кода операции, остальные — для хранения кодов адресов операндов.



2.3. На каких принципах построены компьютеры

В основу построения подавляющего большинства компьютеров положены следующие общие принципы, сформулированные в 1945 г. американским ученым Джоном фон Нейманом.

Принцип программного управления. Из него следует, что программа состоит из набора команд, которые выполняются процессором автоматически друг за другом в определенной последовательности.

Выборка программы из памяти осуществляется с помощью счетчика команд. Этот регистр процессора последовательно увеличивает хранимый в нем адрес очередной команды на длину команды.

А так как команды программы расположены в памяти друг за другом, то тем самым организуется выборка цепочки команд из последовательно расположенных ячеек памяти.

Если же нужно после выполнения команды перейти не к следующей, а к какой-то другой, то используются команды условного или безусловного переходов, которые заносят в счетчик команд номер ячейки памяти, содержащей следующую команду. Выборка команд из памяти прекращается после достижения и выполнения команды "стоп".

Таким образом, процессор исполняет программу автоматически, без вмешательства человека.

Принцип однородности памяти. Программы и данные хранятся в одной и той же памяти. Поэтому компьютер не различает, что хранится в данной ячейке памяти — число, текст или команда. Над командами можно выполнять такие же действия, как и над данными. Это открывает целый ряд возможностей. Например, программа в процессе своего выполнения также может подвергаться переработке, что позволяет задавать в самой программе правила получения некоторых ее частей (так в программе организуется выполнение циклов и подпрограмм). Более того, команды одной программы могут быть получены как результаты исполнения другой программы. На этом принципе основаны методы трансляции — перевода текста программы с языка программирования высокого уровня на язык конкретной машины.

Принцип адресности. Структурно основная память состоит из перенумерованных ячеек; процессору в произвольный момент времени доступна любая ячейка. Отсюда следует возможность давать имена областям памяти, так, чтобы к запомненным в них значениям можно было впоследствии обращаться или менять их в процессе выполнения программ с использованием присвоенных имен.

Компьютеры, построенные на этих принципах, относятся к типу фон-неймановских. Но существуют компьютеры, принципиально отличающиеся от фон-неймановских. Для них, например, может не выполняться принцип программного управления, т.е. они могут работать без “счетчика команд”, указывающего текущую выполняемую команду программы. Для обращения к какой-либо переменной, хранящейся в памяти, этим компьютерам не обязательно давать ей имя. Такие компьютеры называются не-фон-неймановскими.



2.4. Что такое команда

Команда — это описание элементарной операции, которую должен выполнить компьютер. В общем случае, команда содержит следующую информацию:

- код выполняемой операции;
- указания по определению операндов (или их адресов);
- указания по размещению получаемого результата.

В зависимости от количества операндов, команды бывают:

- одноадресные;
- двухадресные;
- трехадресные;
- переменнаяадресные.

Команды хранятся в ячейках памяти в двоичном коде.

В современных компьютерах длина команд переменная (обычно от двух до четырех байтов), а способы указания адресов переменных весьма разнообразные. В адресной части команды может быть указан, например:

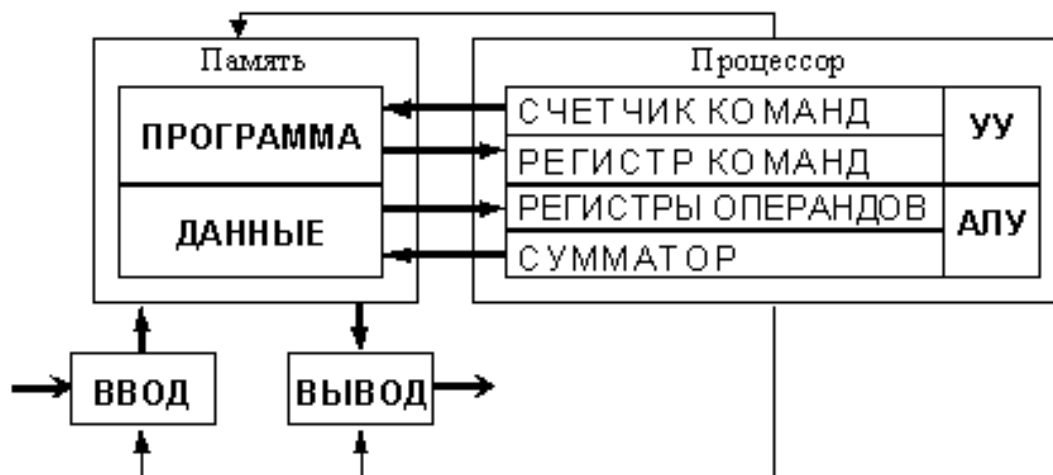
- ⊕ сам операнд (число или символ);
- ⊕ адрес операнда (номер байта, начиная с которого расположен операнд);
- ⊕ адрес адреса операнда (номер байта, начиная с которого расположен адрес операнда), и

др.



2.5. Как выполняется команда

Выполнение команды можно проследить по схеме:



2.2 Общая схема компьютера

Как правило, этот процесс разбивается на следующие этапы:

- из ячейки памяти, адрес которой хранится в счетчике команд, выбирается очередная команда;

- содержимое счетчика команд при этом увеличивается на длину команды;
- выбранная команда передается в устройство управления на регистр команд;
- устройство управления расшифровывает адресное поле команды;
- по сигналам УУ операнды считываются из памяти и записываются в АЛУ на специальные регистры операндов;
- УУ расшифровывает код операции и выдает в АЛУ сигнал выполнить соответствующую операцию над данными;
- результат операции либо остается в процессоре, либо отправляется в память, если в команде был указан адрес результата;
- все предыдущие этапы повторяются до достижения команды “стоп”.



2.6. Что такое архитектура и структура компьютера

При рассмотрении компьютерных устройств, принято различать их архитектуру и структуру.

Архитектурой компьютера называется его описание на некотором общем уровне, включающее описание пользовательских возможностей программирования, системы команд, системы адресации, организации памяти и т.д. Архитектура определяет принципы действия, информационные связи и взаимное соединение основных логических узлов компьютера: процессора, оперативного ЗУ, внешних ЗУ и периферийных устройств. Общность архитектуры разных компьютеров обеспечивает их совместимость с точки зрения пользователя.

Структура компьютера — это совокупность его функциональных элементов и связей между ними. Элементами могут быть самые различные устройства — от основных логических узлов компьютера до простейших схем. Структура компьютера графически представляется в виде структурных схем, с помощью которых можно дать описание компьютера на любом уровне детализации.

Наиболее распространены следующие архитектурные решения.

Классическая архитектура (архитектура фон Неймана) — одно арифметико-логическое устройство (АЛУ), через которое проходит поток данных, и одно устройство управления (УУ), через которое проходит поток команд — программа (рис.2.1). Это однопроцессорный компьютер. К этому типу архитектуры относится и архитектура персонального компьютера с общей шиной. Все функциональные блоки здесь связаны между собой общей шиной, называемой также системной магистралью.

Физически магистраль представляет собой многопроводную линию с гнездами для подключения электронных схем. Совокупность проводов магистрали разделяется на отдельные группы: шину адреса, шину данных и шину управления.

Периферийные устройства (принтер и др.) подключаются к аппаратуре компьютера через специальные контроллеры — устройства управления периферийными устройствами.

Контроллер — устройство, которое связывает периферийное оборудование или каналы связи с центральным процессором, освобождая процессор от непосредственного управления функционированием данного оборудования.

Многопроцессорная архитектура. Наличие в компьютере нескольких процессоров означает, что параллельно может быть организовано много потоков данных и много потоков команд. Таким образом, параллельно могут выполняться несколько фрагментов одной задачи. Структура такой машины, имеющей общую оперативную память и несколько процессоров, представлена на рис. 2.3.

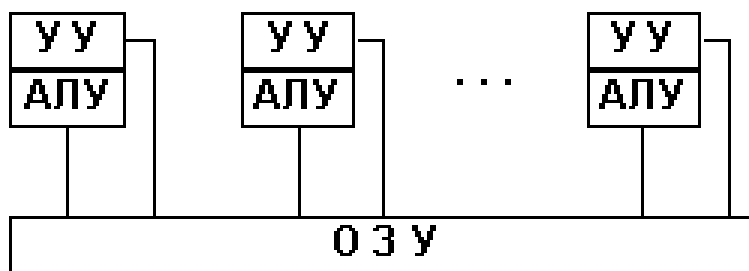


Рис. 2.3. Архитектура многопроцессорного компьютера

Многомашинная вычислительная система. Здесь несколько процессоров, входящих в вычислительную систему, не имеют общей оперативной памяти, а имеют каждый свою (локальную). Каждый компьютер в многомашинной системе имеет классическую архитектуру, и такая система применяется достаточно широко. Однако эффект от применения такой вычислительной системы может быть получен только при решении задач, имеющих очень специальную структуру: она должна разбиваться настолько слабо связанных подзадач, сколько компьютеров в системе.

Архитектура с параллельными процессорами. Здесь несколько АЛУ работают под управлением одного УУ. Это означает, что множество данных может обрабатываться по одной программе — то есть по одному потоку команд. Высокое быстродействие такой архитектуры можно получить только на задачах, в которых одинаковые вычислительные операции выполняются одновременно на различных однотипных наборах данных. Структура таких компьютеров представлена на рис. 2.4.

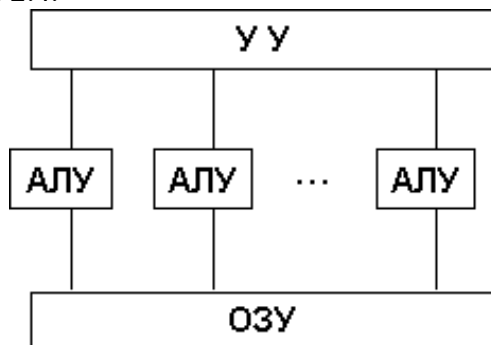


Рис. 2.4. Архитектура с параллельным процессором

В современных машинах часто присутствуют элементы различных типов архитектурных решений. Существуют и такие архитектурные решения, которые радикально отличаются от рассмотренных выше.



2.7. Что такое центральный процессор.

Центральный процессор (CPU от англ. Central Processing Unit) — это основной рабочий компонент компьютера, который выполняет арифметические и логические операции, заданные программой, управляет вычислительным процессом и координирует работу всех устройств компьютера.

Центральный процессор в общем случае содержит в себе:

- арифметико-логическое устройство;
- шины данных и шины адресов;
- регистры;
- счетчики команд;
- кэш — очень быструю память малого объема (от 8 до 512 Кбайт);
- математический сопроцессор чисел с плавающей точкой.

Современные процессоры выполняются в виде микропроцессоров. Физически микропроцессор представляет собой интегральную схему — тонкую пластинку кристаллического кремния прямоугольной формы площадью всего несколько квадратных миллиметров, на которой размещены схемы, реализующие все функции процессора. Кристалл-пластинка обычно помещается в пластмассовый или керамический плоский корпус и соединяется золотыми проводками с металлическими штырьками, чтобы его можно было присоединить к системной плате компьютера.

Микропроцессор Intel Pentium 4 — наиболее совершенный и мощный процессор выпуска 2001 г. с тактовой частотой до 2 Гигагерц. Он предназначен для работы приложений, требующих высокой производительности процессора, таких, как передача видео и звука по Интернет, создание видеоматериалов, распознавание речи, обработка трехмерной графики, игры.

В вычислительной системе может быть несколько параллельно работающих процессоров; такие системы называются многопроцессорными.



2.8. Как устроена память

Память компьютера построена из двоичных запоминающих элементов — битов, объединенных в группы по 8 битов, которые называются байтами. (Единицы измерения памяти совпадают с единицами измерения информации). Все байты пронумерованы. Номер байта называется его адресом.

Байты могут объединяться в ячейки, которые называются также словами. Для каждого компьютера характерна определенная длина слова — два, четыре или восемь байтов. Это не исключает использования ячеек памяти другой длины (например, полуслово, двойное слово). Как правило, в одном машинном слове может быть представлено либо одно целое число, либо одна команда. Однако, допускаются переменные форматы представления информации. Разбиение памяти на слова для четырехбайтовых компьютеров представлено в таблице:

Байт 0	Байт 1	Байт 2	Байт 3	Байт 4	Байт 5	Байт 6	Байт 7
ПОЛУСЛОВО		ПОЛУСЛОВО		ПОЛУСЛОВО		ПОЛУСЛОВО	
СЛОВО				СЛОВО			
ДВОЙНОЕ СЛОВО							

Широко используются и более крупные производные единицы объема памяти: Килобайт, Мегабайт, Гигабайт, а также Терабайт и Петабайт.

Современные компьютеры имеют много разнообразных запоминающих устройств, которые сильно отличаются между собой по назначению, временным характеристикам, объёму хранимой информации и стоимости хранения одинакового объёма информации. Различают два основных вида памяти — внутреннюю и внешнюю.



2.9. Какие устройства образуют внутреннюю память

В состав внутренней памяти входят оперативная память, кэш-память и специальная память.

1. Оперативная память

Оперативная память (ОЗУ англ. RAM, Random Access Memory — память с произвольным доступом) — это быстрое запоминающее устройство не очень большого объёма, непосредственно связанное с процессором и предназначенное для записи, считывания и хранения выполняемых программ и данных, обрабатываемых этими программами.

Оперативная память используется только для временного хранения данных и программ, так как, в момент включения машины, все, что находилось в ОЗУ, пропадает. Доступ к элементам оперативной памяти прямой — это означает, что каждый байт памяти имеет свой индивидуальный адрес.

Объем ОЗУ обычно составляет от 32 до 512 Мбайт. Для несложных административных задач бывает достаточно и 32 Мбайт ОЗУ, но сложные задачи компьютерного дизайна могут потребовать от 512 Мбайт до 2 Гбайт ОЗУ.

Обычно ОЗУ исполняется из интегральных микросхем памяти SDRAM (синхронное динамическое ОЗУ). Каждый информационный бит в SDRAM запоминается в виде электрического заряда крохотного конденсатора, образованного в структуре полупроводникового кристалла. Из-за токов утечки такие конденсаторы быстро разряжаются, и их периодически (примерно каждые 2 миллисекунды) подзаряжают специальные устройства. Этот процесс называется регенерацией памяти (Refresh Memory). Микросхемы SDRAM имеют ёмкость 16 — 256 Мбит и более. Они устанавливаются в корпуса и собираются в модули памяти.

Большинство современных компьютеров комплектуются модулями типа DIMM (Dual-In-line Memory Module — модуль памяти с двухрядным расположением микросхем). В компьютерных системах на самых современных процессорах используются высокоскоростные модули Rambus DRAM (RIMM) и DDR DRAM.

Модули памяти характеризуются такими параметрами, как объем — (16, 32, 64, 128, 256 или 512 Мбайт), число микросхем, паспортная частота (100 или 133 МГц), время доступа к данным (6 или 7 наносекунд) и число контактов (72, 168 или 184).

2. Кэш-память

Кэш (англ. - cache), или сверхоперативная память — очень быстрое ЗУ небольшого объема, которое используется при обмене данными между микропроцессором и оперативной памятью для компенсации разницы в скорости обработки информации процессором и несколько менее быстросрающей оперативной памятью.

Кэш-памятью управляет специальное устройство — контроллер, который, анализируя выполняемую программу, пытается предвидеть, какие данные и команды вероятнее всего понадобятся в ближайшее время процессору, и подкачивает их в кэш-память. При этом возможны как "попадания", так и "промахи". В случае попадания, то есть, если в кэш подкачаны нужные данные, извлечение их из памяти происходит без задержки. Если же требуемая информация в кэше отсутствует, то процессор считывает её непосредственно из оперативной памяти. Соотношение числа попаданий и промахов определяет эффективность кэширования.

Кэш-память реализуется на микросхемах статической памяти SRAM (Static RAM), более быстросрающих, дорогих и малоёмких, чем DRAM (SDRAM). Современные микропроцессоры имеют встроенную кэш-память, так называемый кэш первого уровня размером 8, 16 или 32 Кбайт. Кроме того, на системной плате компьютера может быть установлен кэш второго уровня ёмкостью 256, 512 Кбайт и выше.

3. Специальная память

К устройствам специальной памяти относятся постоянная память (ROM), перепрограммируемая постоянная память (Flash Memory), память CMOS RAM, питаемая от батарейки, видеопамять и некоторые другие виды памяти.

Постоянная память (ПЗУ, англ. ROM, Read Only Memory — память только для чтения) — энергонезависимая память, используется для хранения данных, которые никогда не требуют изменения. Содержание памяти специальным образом "зашивается" в устройстве при его изготовлении для постоянного хранения. Из ПЗУ можно только читать.

Перепрограммируемая постоянная память (Flash Memory) — энергонезависимая память, допускающая многократную перезапись своего содержимого с дискеты.

Прежде всего, в постоянную память записывают программу управления работой самого процессора. В ПЗУ находятся программы управления дисплеем, клавиатурой, принтером, внешней памятью, программы запуска и остановки компьютера, тестирования устройств.

Важнейшая микросхема постоянной или Flash-памяти — модуль BIOS. Роль BIOS двоякая: с одной стороны это неотъемлемый элемент аппаратуры, а с другой стороны — важный модуль любой операционной системы.

BIOS (Basic Input/Output System — базовая система ввода-вывода) — совокупность программ, предназначенных для автоматического тестирования устройств после включения питания компьютера и загрузки операционной системы в оперативную память.

Разновидность постоянного ЗУ — CMOS RAM.

CMOS RAM — это память с невысоким быстросраждением и минимальным энергопотреблением от батарейки. Используется для хранения информации о конфигурации и составе оборудования компьютера, а также о режимах его работы.

Интегральные схемы BIOS и CMOS

Содержимое CMOS изменяется специальной программой Setup, находящейся в BIOS (англ. Set-up — устанавливать, читается "сетап").

Для хранения графической информации используется видеопамять.

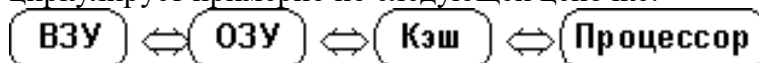
Видеопамять (VRAM) — разновидность оперативного ЗУ, в котором хранятся закодированные изображения. Это ЗУ организовано так, что его содержимое доступно сразу двум устройствам — процессору и дисплею. Поэтому изображение на экране меняется одновременно с обновлением видеоданных в памяти.





2.10. Какие устройства образуют внешнюю память

Внешняя память (ВЗУ) предназначена для длительного хранения программ и данных, и целостность её содержимого не зависит от того, включен или выключен компьютер. В отличие от оперативной памяти, внешняя память не имеет прямой связи с процессором. Информация от ВЗУ к процессору и наоборот циркулирует примерно по следующей цепочке:



В состав внешней памяти компьютера входят:

- накопители на жёстких магнитных дисках;
- накопители на гибких магнитных дисках;
- накопители на компакт-дисках;
- накопители на магнито-оптических компакт-дисках;
- накопители на магнитной ленте (стримеры) и др.

1. Накопители на жёстких магнитных дисках

Если гибкие диски — это средство переноса данных между компьютерами, то жесткий диск — информационный склад компьютера.

Накопитель на жёстких магнитных дисках (англ. HDD — Hard Disk Drive) или винчестерский накопитель — это наиболее массовое запоминающее устройство большой ёмкости, в котором носителями информации являются круглые алюминиевые пластины — платтеры, обе поверхности которых покрыты слоем магнитного материала. Используется для постоянного хранения информации — программ и данных.

Как и у дискеты, рабочие поверхности платтеров разделены на кольцевые концентрические дорожки, а дорожки — на секторы. Головки считывания-записи вместе с их несущей конструкцией и дисками заключены в герметически закрытый корпус, называемый модулем данных. При установке модуля данных на дисковод он автоматически соединяется с системой, подкачивающей очищенный охлажденный воздух. Поверхность платтера имеет магнитное покрытие толщиной всего лишь в 1,1 мкм, а также слой смазки для предохранения головки от повреждения при опускании и подъёме на ходу. При вращении платтера над ним образуется воздушный слой, который обеспечивает воздушную подушку для зависания головки на высоте 0,5 мкм над поверхностью диска.

Винчестерские накопители имеют очень большую ёмкость: от 10 до 100 Гбайт. У современных моделей скорость вращения шпинделя (вращающего вала) обычно составляет 7200 об/мин, среднее время поиска данных 9 мс, средняя скорость передачи данных до 60 Мбайт/с. В отличие от дискеты, жесткий диск вращается непрерывно. Все современные накопители снабжаются встроенным кэшем (обычно 2 Мбайта), который существенно повышает их производительность. Винчестерский накопитель связан с процессором через контроллер жесткого диска.

3. Накопители на компакт-дисках

Здесь носителем информации является *CD-ROM* (Compact Disk Read-Only Memory - компакт диск, из которого можно только читать). *CD-ROM RW* — можно записывать и перезаписывать.

CD-ROM представляет собой прозрачный полимерный диск диаметром 12 см и толщиной 1,2 мм, на одну сторону которого напылен светоотражающий слой алюминия, защищенный от повреждений слоем прозрачного лака. Толщина напыления составляет несколько десятитысячных долей миллиметра.

Информация на диске представляется в виде последовательности впадин (углублений в диске) и выступов (их уровень соответствует поверхности диска), расположенных на спиральной дорожке, выходящей из области вблизи оси диска. На каждом дюйме (2,54 см) по радиусу диска размещается 16 тысяч витков спиральной дорожки. Для сравнения — на поверхности жесткого диска на дюйме по радиусу помещается лишь несколько сотен дорожек. Емкость CD достигает 780 Мбайт. Информация наносится на диск при его изготовлении и не может быть изменена.

CD-ROM обладают высокой удельной информационной емкостью, что позволяет создавать на их основе справочные системы и учебные комплексы с большой иллюстративной базой. Один CD по информационной емкости равен почти 500 дискетам. Считывание информации с CD-ROM происходит с достаточно высокой скоростью, хотя и заметно меньшей, чем скорость работы накопителей на жестком диске. CD-ROM просты и удобны в работе, имеют низкую удельную стоимость хранения данных, практически не изнашиваются, не могут быть поражены вирусами, с них невозможно случайно стереть информацию.

В отличие от магнитных дисков, компакт-диски имеют не множество кольцевых дорожек, а одну — спиральную, как у грампластинок. В связи с этим, угловая скорость вращения диска не постоянна. Она линейно уменьшается в процессе продвижения читающей лазерной головки к краю диска.

Для работы с CD-ROM нужно подключить к компьютеру накопитель CD-ROM, преобразующий последовательность углублений и выступов на поверхности CD-ROM в последовательность двоичных сигналов. Для этого используется считывающая головка с микролазером и светодиодом. Глубина впадин на поверхности диска равна четверти длины волны лазерного света. Если в двух последовательных тактах считывания информации луч света лазерной головки переходит с выступа на дно впадины или обратно, разность длин путей света в этих тактах меняется на полуволну, что вызывает усиление или ослабление совместно попадающих на светодиод прямого и отраженного от диска света. Если в последовательных тактах считывания длина пути света не меняется, то и состояние светодиода не меняется. В результате ток через светодиод образует последовательность двоичных электрических сигналов, соответствующих сочетанию впадин и выступов на дорожке.



Различная длина оптического пути луча света в двух последовательных тактах считывания информации соответствует двоичным единицам. Одинаковая длина соответствует двоичным нулям. Сегодня почти все персональные компьютеры имеют накопитель CD-ROM. Но многие мультимедийные интерактивные программы слишком велики, чтобы поместиться на одном CD. На смену технологии CD-ROM стремительно идет технология цифровых видеодисков DVD. Эти диски имеют тот же размер, что и обычные CD, но вмещают до 17 Гбайт данных, т.е. по объему заменяют 20 стандартных дисков CD-ROM. На таких дисках выпускаются мультимедийные игры и интерактивные видеофильмы отличного качества, позволяющие зрителю просматривать эпизоды под разными углами камеры, выбирать различные варианты окончания картины, знакомиться с биографиями снявшихся актеров, наслаждаться великолепным качеством звука.

Архитектура, назначение и функции операционных систем

1.1. Понятие операционной системы. Виртуальные машины

Современный компьютер — сложнейшая аппаратно-программная система. Написание программ для компьютера, их отладка и последующее выполнение представляет собой сложную трудоемкую задачу. Основная причина этого — огромная разница между тем, что удобно для людей, и тем, что удобно для компьютеров. Компьютер понимает только свой, машинный язык (назовем его Я0), а для человека наиболее удобен разговорный или хотя бы язык описания алгоритмов — алгоритмический язык. Проблему можно решить двумя способами. Оба способа связаны с разработкой команд, которые были бы более удобны для человека, чем встроенные

машинные команды компьютера. Эти новые команды в совокупности формируют некоторый язык, который назовем Я1.

Упомянутые два способа решения проблемы различаются тем, каким образом компьютер будет выполнять программы, написанные на языке Я1. Первый способ – замена каждой команды языка Я1 на эквивалентный набор команд в языке Я0. В этом случае компьютер выполняет новую программу, написанную на языке Я0, вместо программы, написанной на языке Я1. Эта технология называется *трансляцией*.

Второй способ – написание программы на языке Я0, которая берет программы, написанные на языке Я1, в качестве входных данных, рассматривает каждую команду по очереди и сразу выполняет эквивалентный набор команд языка Я0. Эта технология не требует составления новой программы на Я0. Она называется *интерпретацией*, а программа, которая осуществляет интерпретацию, называется *интерпретатором*.

В подобной ситуации проще представить себе существование гипотетического компьютера или *виртуальной* машины, для которой машинным языком является язык Я1, чем думать о трансляции и интерпретации. Назовем такую виртуальную машину М1, а виртуальную машину с языком Я0 – М0. Для виртуальных машин можно будет писать программы, как будто они (машины) действительно существуют.

Очевидно, можно пойти дальше – создать еще набор команд, который в большей степени ориентирован на человека и в меньшей степени на компьютер, чем Я1. Этот набор формирует язык Я2 и, соответственно, виртуальную машину М2. Так можно продолжать до тех пор, пока не дойдем до подходящего нам языка уровня n.

Большинство современных компьютеров состоит из двух и более уровней. Уровень 0 – аппаратное обеспечение машины. Электронные схемы этого уровня выполняют программы, написанные на языке уровня 1.

Следующий уровень – *микроархитектурный* уровень. На этом уровне можно видеть совокупности 8 или 32 (иногда и больше) регистров, которые формируют локальную память и АЛУ (арифметико-логическое устройство). Регистры вместе с АЛУ формируют тракт данных, по которому поступают данные. Основная операция этого тракта заключается в следующем. Выбирается один или два регистра, АЛУ производит над ними какую-то операцию, а результат помещается в один из этих регистров. На некоторых машинах работа тракта контролируется особой программой, которая называется микропрограммой. В других машинах такой контроль выполняется аппаратным обеспечением.

Следующий (второй) уровень составляет уровень *архитектуры системы команд*. Команды используют регистры и другие возможности аппаратуры. Команды формируют уровень ISA (Instruction Set Architecture), называемый машинным языком. Обычно машинный язык содержит от 50 до 300 команд, служащих преимущественно для перемещения данных по компьютеру, выполнения арифметических операций и сравнения величин.

Следующий (третий) уровень обычно – гибридный. Большинство команд в его языке есть также и на уровне архитектуры системы команд. У этого уровня есть некоторые дополнительные особенности: набор новых команд, другая организация памяти, способность выполнять две и более программы одновременно и некоторые другие. С течением времени набор таких команд существенно расширился. В нем появились так называемые макросы операционной системы или вызовы супервизора, называемые теперь системными вызовами.

Новые средства, появившиеся на третьем уровне, выполняются интерпретатором, который работает на втором уровне. Этот интерпретатор был когда-то назван *операционной системой*. Команды третьего уровня, идентичные командам второго уровня, выполняются микропрограммой или аппаратным обеспечением, но не операционной системой. Иными словами, одна часть команд третьего уровня интерпретируется операционной системой, а другая часть – микропрограммой. Вот почему этот уровень операционной системы считается гибридным.

Операционная система была создана для того, чтобы автоматизировать работу оператора и скрыть от пользователя сложности общения с аппаратурой, предоставив ему более удобную систему команд. Нижние три уровня (с нулевого по второй) конструируются не для того, чтобы с ними работал обычный программист. Они изначально предназначены для работы интерпретаторов и трансляторов, поддерживающих более высокие уровни. Эти трансляторы и интерпретаторы

составляются системными программистами, которые специализируются на разработке и построении новых виртуальных машин.

Над операционной системой (ОС) расположены остальные системные программы. Здесь находятся интерпретатор команд (оболочка), компиляторы, редакторы и т.д. Подобные программы не являются частью ОС (иногда оболочку пользователи считают операционной системой). Под операционной системой обычно понимается то программное обеспечение, которое запускается в *режиме ядра* или, как еще его называют, *режиме супервизора*. Она защищена от вмешательства пользователя с помощью специальных аппаратных средств.

Четвертый уровень представляет собой символическую форму одного из языков низкого уровня (обычно ассемблер). На этом уровне можно писать программы в приемлемой для человека форме. Эти программы сначала транслируются на язык уровня 1, 2 или 3, а затем интерпретируются соответствующей виртуальной или фактически существующей (физической) машиной.

Уровни с пятого и выше предназначены для прикладных программистов, решающих конкретные задачи на языках высокого уровня (C, C++, C#, VBA и др.). Компиляторы и редакторы этих уровней *запускаются в пользовательском режиме*. На еще более высоких уровнях располагаются прикладные программы пользователей.

Большинство пользователей компьютеров имеют опыт общения с операционной системой, по крайней мере, в той степени, чтобы эффективно выполнять свои текущие задачи. Однако они испытывают затруднения при попытке дать определение операционной системе. В известной степени проблема связана с тем, что операционные системы выполняют две основные, но практически не связанные между собой функции: расширение возможностей компьютера и управление его ресурсами.

С точки зрения пользователя ОС выполняет функцию расширенной машины или виртуальной машины, в которой легче программировать и легче работать, чем непосредственно с аппаратным обеспечением, составляющим реальный компьютер. Операционная система не только устраняет необходимость работы непосредственно с дисками и предоставляет простой, ориентированный на работу с файлами интерфейс, но и скрывает множество неприятной работы с прерываниями, счетчиками времени, организацией памяти и другими компонентами низкого уровня.

Однако концепция, рассматривающая операционную систему прежде всего как удобный интерфейс пользователя, – это взгляд сверху вниз. Альтернативный взгляд, снизу вверх, дает представление об операционной системе как о механизме, присутствующем в компьютере для управления всеми компонентами этой сложнейшей системы. В соответствии с этим подходом работа операционной системы заключается в обеспечении организованного и контролируемого распределения процессоров, памяти, дисков, принтеров, устройств ввода-вывода, датчиков времени и т.п. между различными программами, конкурирующими за право их использовать.

1.2. Операционная система, среда и операционная оболочка

Операционные системы (ОС) в современном их понимании (их назначении и сущности) появились значительно позже первых компьютеров (правда, по всей видимости, и исчезнут в этой сущности в компьютерах будущего). Почему и когда появились ОС? Считается что первая цифровая вычислительная машина ENIAC (Electronic Numerical Integrator and Computer) была создана в 1946 году по проекту "Проект PX" Министерства обороны США. На реализацию проекта затрачено 500 тыс. долларов. Компьютер содержал 18000 электронных ламп, массу всякой электроники, включал в себя 12 десятиразрядных сумматоров, а для ускорения некоторых арифметических операций имел умножитель и "делитель-извлекающий" квадратного корня. Программирование сводилось к связыванию различных блоков проводами. Конечно, никакого программного обеспечения и тем более операционных систем тогда еще не существовало.

Интенсивное создание различных моделей ЭВМ относится к началу 50-х годов прошлого века. В эти годы одни и те же группы людей участвовали и в проектировании, и в создании, и в программировании, и в эксплуатации ЭВМ. Программирование осуществлялось исключительно на машинном языке (а затем на ассемблере), не было никакого системного программного

обеспечения, кроме библиотек математических и служебных подпрограмм. Операционные системы еще не появились, а все задачи организации вычислительного процесса решались вручную каждым программистом с примитивного пульта управления ЭВМ.

С появлением полупроводниковых элементов вычислительные возможности компьютеров существенно выросли. Наряду с этим заметно прогрессируют достижения в области автоматизации программирования и организации вычислительных работ. Появились алгоритмические языки (алгол, фортран, кобол) и системное программное обеспечение (трансляторы, редакторы связи, загрузчики и др.). Выполнение программ усложнилось и включало в себя следующие основные действия:

- загрузка нужного транслятора;
- запуск транслятора и получение программы в машинных кодах;
- связывание программы с библиотечными подпрограммами;
- загрузка программы в оперативную память;
- запуск программы;
- вывод результатов работы программы на печатающее или другое периферийное устройство.

Для организации эффективной загрузки всех средств компьютера в штаты вычислительных центров ввели должности специально обученных операторов, профессионально выполнявших работу по организации вычислительного процесса для всех пользователей этого центра. Однако, как бы ни был подготовлен оператор, ему тяжело состязаться в производительности с работой устройств компьютера. И поэтому большую часть времени дорогостоящий процессор простаивал, а следовательно, использование компьютеров не было эффективным.

С целью исключения простоев были предприняты попытки разработки специальных программ – мониторов, прообразов первых операционных систем, которые осуществляли автоматический переход от задания к заданию. Считается, что первую операционную систему создала в 1952 году для своих компьютеров IBM-701 исследовательская лаборатория фирмы General Motors [9]. В 1955 году эта фирма и North American Aviation совместно разработали ОС для компьютера IBM-704.

В конце 50-х годов прошлого века ведущие фирмы изготовители поставляли операционные системы со следующими характеристиками:

- пакетная обработка одного потока задач;
- наличие стандартных программ ввода-вывода;
- возможности автоматического перехода от программы к программе;
- средства восстановления после ошибок, обеспечивающие автоматическую "очистку" компьютера в случае аварийного завершения очередной задачи и позволяющие запускать следующую задачу при минимальном вмешательстве оператора;
- языки управления заданиями, предоставляющие пользователям возможность описывать свои задания и ресурсы, требуемые для их выполнения.

Пакет представляет собой набор (колоду) перфокарт, организованную специальным образом (задание, программы, данные). Для ускорения работы он мог переноситься на магнитную ленту или диск. Это позволяло сократить простой дорогой аппаратуры. Надо сказать, что в настоящее время в связи с прогрессом микроэлектронных технологий и методологий программирования значительно снизилась стоимость аппаратных и программных средств компьютерной техники. Поэтому сейчас основное внимание уделяется тому, чтобы сделать работу пользователей и программистов более эффективной, поскольку затраты труда квалифицированных специалистов сейчас представляют собой гораздо большую долю общей стоимости вычислительных систем, чем аппаратные и программные средства компьютеров.

Расположение операционной системы в иерархической структуре программного и аппаратного обеспечения компьютера можно представить, как показано на рис. 1.1.



Рис. 1.1. Иерархическая структура программно-аппаратных средств компьютера

Самый нижний уровень содержит различные устройства компьютера, состоящие из микросхем, проводников, источников питания, электронно-лучевых трубок и т.п. Этот уровень можно разделить на подуровни, например контроллеры устройств, а затем сами устройства. Возможно деление и на большее число уровней. Выше расположен микроархитектурный уровень, на котором физические устройства рассматриваются как отдельные функциональные единицы.

На микроархитектурном уровне находятся внутренние регистры центрального процессора (их может быть несколько) и арифметико-логические устройства со средствами управления ими. На этом уровне реализуется выполнение машинных команд. В процессе выполнения команд используются регистры процессора и устройств, а также другие возможности аппаратуры. Команды, видимые для работающего на ассемблере программиста, формируют уровень ISA (Instruction Set Architecture – архитектура системы команд), часто называемый машинным языком.

Операционная система предназначена для того, чтобы скрыть все эти сложности. Конечный пользователь обычно не интересуется деталями устройства аппаратного обеспечения компьютера. Компьютер ему видится как набор приложений. Приложение может быть написано программистом на каком-либо языке программирования. Для упрощения этой работы программист использует набор системных программ, некоторые из которых называются утилитами. С их помощью реализуются часто используемые функции, которые помогают работать с файлами, управлять устройствами ввода-вывода и т.п. Программист применяет эти средства при разработке программ, а приложения во время выполнения обращаются к утилитам для выполнения определенных функций. Наиболее важной из системных программ является операционная система, которая освобождает программиста от необходимости глубокого знания устройства компьютера и представляет ему удобный интерфейс для его использования. Операционная система выступает в роли посредника, облегчая программисту, пользователям и программным приложениям доступ к различным службам и возможностям компьютера.

Таким образом, *операционная система* – это набор программ, контролирующих работу прикладных программ и системных приложений и исполняющих роль интерфейса между пользователями, программистами, прикладными программами, системными приложениями и аппаратным обеспечением компьютера.

Образно можно сказать, что аппаратура компьютера предоставляет "сырую" вычислительную мощность, а задача операционной системы заключается в том, чтобы сделать использование этой вычислительной мощности доступным и по возможности удобным для пользователя. Программист может не знать детали управления конкретными ресурсами (например, диском) компьютера и должен обращаться к операционной системе с соответствующими вызовами, чтобы получить от нее необходимые сервисы и функции. Этот набор сервисов и функций и представляет собой операционную среду, в которой выполняются прикладные программы.

Таким образом, *операционная среда* – это программная среда, образуемая операционной системой, определяющая интерфейс прикладного программирования (API) как множество

системных функций и сервисов (системных вызовов), которые предоставляются прикладным программам. Операционная среда может включать несколько интерфейсов прикладного программирования. Кроме основной операционной среды, называемой естественной (native), могут быть организованы путем эмуляции (моделирования) дополнительные программные среды, позволяющие выполнять приложения, которые рассчитаны на другие операционные системы и даже другие компьютеры.

Еще одно важное понятие, связанное с операционной системой, относится к реализации пользовательских интерфейсов. Как правило, любая операционная система обеспечивает удобную работу пользователя за счет средств пользовательского интерфейса. Эти средства могут быть неотъемлемой частью операционной среды (например, графический интерфейс Windows или текстовый интерфейс командной строки MS DOS), а могут быть реализованы отдельной системной программой – оболочкой операционной системы (например, Norton Commander для MS DOS). В общем случае под *оболочкой операционной системы* понимается часть операционной среды, определяющая интерфейс пользователя, его реализацию (текстовый, графический и т.п.), командные и сервисные возможности пользователя по управлению прикладными программами и компьютером.

Перейдем к рассмотрению эволюции операционных систем.

1.3. Эволюция операционных систем

Рассматривая эволюцию ОС, следует иметь в виду, что разница во времени реализации некоторых принципов организации отдельных операционных систем до их общего признания, а также терминологическая неопределенность не позволяют дать точную хронологию развития ОС. Однако сейчас уже достаточно точно можно определить основные вехи на пути эволюции операционных систем.

Существуют также различные подходы к определению поколений ОС. Известно деление ОС на поколения в соответствии с поколениями вычислительных машин и систем. Такое деление нельзя считать полностью удовлетворительным, так как развитие методов организации ОС в рамках одного поколения ЭВМ, как показал опыт их создания, происходит в достаточно широком диапазоне. Другая точка зрения не связывает поколение ОС с соответствующими поколениями ЭВМ. Так, например, известно определение поколений ОС по уровням входного языка ЭВМ, режимам использования центральных процессоров, формам эксплуатации систем и т.п.

Видимо, наиболее целесообразным следует считать выделение этапов развития ОС в рамках отдельных поколений ЭВМ и ВС.

Первым этапом развития системного программного обеспечения можно считать использование библиотечных программ, стандартных и служебных подпрограмм и макрокоманд. Концепция библиотек подпрограмм является наиболее ранней и восходит к 1949 году. С появлением библиотек получили развитие автоматические средства их сопровождения – программы-загрузчики и редакторы связей. Эти средства применялись в ЭВМ первого поколения, когда операционных систем как таковых еще не существовало.

Стремление устранить несоответствие между производительностью процессоров и скоростью работы электромеханических устройств ввода-вывода, с одной стороны, и использование достаточно быстродействующих накопителей на магнитных лентах и барабанах (НМЛ и НМБ), а затем на магнитных дисках (НМД), с другой стороны, привело к необходимости решения задач буферизации и блокирования-деблокирования данных. Возникли специальные программы методов доступа, которые вносились в объекты модулей редакторов связей (впоследствии стали использоваться принципы полибуферизации). Для поддержания работоспособности и облегчения процессов эксплуатации машин создавались диагностические программы. Таким образом было создано базовое системное программное обеспечение.

С улучшением характеристик ЭВМ и ростом их производительности стало ясно, что существующего базового программного обеспечения (ПО) недостаточно. Появились операционные системы ранней пакетной обработки – мониторы. В рамках системы пакетной обработки во время выполнения любой работы в пакете (трансляция, сборка, выполнение готовой программы) никакая часть системного ПО не находилась в оперативной памяти, так как вся

память предоставлялась текущей работе. Затем появились мониторные системы, в которых оперативная память делилась на три области: фиксированная область мониторной системы, область пользователя и область общей памяти (для хранения данных, которыми могут обмениваться объектные модули).

Началось интенсивное развитие методов управления данными, возникала такая важная функция ОС, как реализация ввода-вывода без участия центрального процесса – так называемый спулинг (от англ. SPOOL – Simultaneous Peripheral Operation on Line).

Появление новых аппаратных разработок (1959-1963 гг.) – систем прерываний, таймеров, каналов – стимулировало дальнейшее развитие ОС. Возникли исполнительные системы, которые представляли собой набор программ для распределения ресурсов ЭВМ, связей с оператором, управления вычислительным процессом и управления вводом-выводом. Такие исполнительные системы позволили реализовать довольно эффективную по тому времени форму эксплуатации вычислительной системы – однопрограммную пакетную обработку. Эти системы давали пользователю такие средства, как контрольные точки, логические таймеры, возможность построения программ оверлейной структуры, обнаружение нарушений программами ограничений, принятых в системе, управление файлами, сбор учетной информации и др.

Однако однопрограммная пакетная обработка с ростом производительности ЭВМ не могла обеспечить экономически приемлемый уровень эксплуатации машин. Решением стало мультипрограммирование – способ организации вычислительного процесса, при котором в памяти компьютера находится несколько программ, попеременно выполняющихся одним процессором, причем для начала или продолжения счета по одной программе не требовалось завершения других. В мультипрограммной среде проблемы распределения ресурсов и защиты стали более острыми и трудноразрешимыми.

Теория построения операционных систем в этот период обогатилась рядом плодотворных идей. Появились различные формы мультипрограммных режимов работы, в том числе разделение времени – режим, обеспечивающий работу многотерминальной системы. Была создана и развита концепция виртуальной памяти, а затем и виртуальных машин. Режим разделения времени позволил пользователю интерактивно взаимодействовать со своими программами, как это было до появления систем пакетной обработки.

Одной из первых ОС, использующих эти новейшие решения, была операционная система MCP (главная управляющая программа), созданная фирмой Burroughs для своих компьютеров B5000 в 1963 году. В этой ОС были реализованы многие концепции и идеи, ставшие впоследствии стандартными для многих операционных систем:

- мультипрограммирование;
- мультипроцессорная обработка;
- виртуальная память;
- возможность отладки программ на исходном языке;
- написание операционной системы на языке высокого уровня.

Известной системой разделения времени того периода стала система CTSS (Compatible Time Sharing System) – совместимая система разделения времени, разработанная в Массачусетском технологическом институте (1963 год) для компьютера IBM-7094. Эта система была использована для разработки в этом же институте совместно с Bell Labs и General Electric системы разделения времени следующего поколения MULTICS (Multiplexed Information And Computing Service). Примечательно, что эта ОС была написана в основном на языке высокого уровня EPL (первая версия языка PL/1 фирма IBM).

Одним из важнейших событий в истории операционных систем считается появление в 1964 году семейства компьютеров под названием System/360 фирмы IBM, а позже – System/370. Это было первой в мире реализацией концепции семейства программного и информационно совместимых компьютеров, ставшей впоследствии стандартной для всех фирм компьютерной отрасли.

Основной формой использования ЭВМ, как в системах разделения времени, так и в системах пакетной обработки, стал многотерминальный режим. При этом не только оператор, но и все пользователи получали возможность формулировать свои задания и управлять их выполнением со своего терминала. Поскольку терминальные комплексы скоро стало возможным

размещать на значительных расстояниях от компьютера (благодаря модемным телефонным соединениям), появились системы удаленного ввода заданий и телеобработки данных. В ОС добавились модули, реализующие протоколы связи.

К этому времени произошло существенное изменение в распределении функций между аппаратными и программными средствами компьютера. Операционная система становится "неотъемлемой частью ЭВМ", как бы продолжением аппаратуры. В процессорах появился привилегированный (Супервизор в OS/360) и пользовательский (Задача в OS/360) режимы работы, мощная система прерываний, защита памяти, специальные регистры для быстрого переключения программ, средства поддержки виртуальной памяти и др.

В начале 70-х годов появились первые сетевые ОС, которые позволили не только рассредоточить пользователей, как в системах телеобработки данных, но и организовать распределенное хранение и обработку данных между компьютерами, соединенных электрическими связями. Известен проект ARPANET МО США. В 1974 году IBM объявила о создании собственной сетевой архитектуры SNA для своих мэйнфреймов, обеспечивающей взаимодействие типа "терминал-терминал", "терминал-компьютер", "компьютер-компьютер". В Европе активно разрабатывалась технология построения сетей с коммутацией пакетов на основе протоколов X.25.

К середине 70-х годов наряду с мэйнфреймами широкое распространение получили мини-компьютеры (PDP-11, Nova, HP). Архитектура мини-компьютеров была значительно проще, многие функции мультипрограммных ОС мэйнфреймов были усечены. Операционные системы мини-ЭВМ стали делать специализированными (RSX-11M – разделение времени, RT-11 – ОС реального времени) и не всегда многопользовательскими.

Важной вехой в истории мини-компьютеров и вообще в истории операционных систем явилось создание ОС UNIX. Написал эту систему Кен Томпсон (Ken Thompson), один из специалистов по компьютерам в BELL Labs, работавший над проектом MULTICS. Собственно, его UNIX – это усеченная однопользовательская версия системы MULTICS. Первоначальное название этой системы – UNICS (UNiplexed Information and Computing Service – примитивная информационная и компьютерная служба). Так в шутку была названа эта система, поскольку MULTICS (MULTiplexed Information and Computing Service) – мультиплексная информационная и компьютерная служба. С середины 70-х годов началось массовое использование ОС UNIX, написанной на 90% на языке C. Широкое распространение C-компиляторов сделало UNIX уникальной переносимой ОС, а поскольку она поставлялась вместе с исходными кодами, она стала первой открытой операционной системой. Гибкость, элегантность, мощные функциональные возможности и открытость позволили ей занять прочные позиции во всех классах компьютеров – от персональных до супер-ЭВМ.

Доступность мини-компьютеров послужила стимулом для создания локальных сетей. В простейших ЛВС компьютеры соединялись через последовательные порты. Первое сетевое приложение для ОС UNIX – программа UUCP (Unix to Unix Copy Program) – появилось в 1976 году.

Дальнейшее развитие сетевых систем со стеком протоколов TCP/IP: в 1983 году он был принят МО США в качестве стандарта и использован в сети ARPANET. В этом же году ARPANET разделилась на MILNET (для военного ведомства США) и новую ARPANET, которую стали называть Internet.

Все восьмидесятые годы характерны появлением все более совершенных версий UNIX: Sun OS, HP-UX, Irix, AIX и др. Для решения проблемы их совместимости были приняты стандарты POSIX и XPG, определяющие интерфейсы этих систем для приложений.

Еще одним знаменательным событием для истории операционных систем было появление в начале 80-х годов персональных компьютеров. Они послужили мощным толчком для распределения локальных сетей, в результате поддержка сетевых функций стала для ОС ПК необходимым условием. Однако и дружелюбный интерфейс, и сетевые функции появились у ОС ПК не сразу.

Наиболее популярной версией ОС раннего этапа развития персональных компьютеров была MS-DOS компании Microsoft – однопрограммная, однопользовательская ОС с интерфейсом командной строки. Многие функции, обеспечивающие удобство работы пользователю, в этой ОС

предоставлялись дополнительными программами – оболочкой Norton Commander, PC Tools и др. Наибольшее влияние на развитие программного обеспечения ПК оказала операционная среда Windows, первая версия которой появилась в 1985 году. Сетевые функции также реализовались с помощью сетевых оболочек и появились в MS-DOS версии 3.1. В это же время появились сетевые продукты Microsoft – MS-NET, а позже – LAN Manager, Windows for Workgroup, а затем и Windows NT.

Другим путем пошла компания Novell: ее продукт NetWare – операционная система со встроенными сетевыми функциями. ОС NetWare распространялась как операционная система для центрального сервера локальной сети и за счет специализации функций файл-сервера обеспечивала высокую скорость удаленного доступа к файлам и повышенную безопасность данных. Однако эта ОС имела специфический программный интерфейс (API), что затрудняло разработку приложений.

В 1987 году появилась первая многозадачная ОС для ПК – OS/2, разработанная Microsoft совместно с IBM. Эта была хорошо продуманная система с виртуальной памятью, графическим интерфейсом и возможностью выполнять DOS-приложения. Для нее были созданы и получили распространение сетевые оболочки LAN Manager (Microsoft) и LAN Server (IBM). Эти оболочки уступали по производительности файловому серверу NetWare и потребляли больше аппаратных ресурсов, но имели важные достоинства. Они позволяли выполнять на сервере любые программы, разработанные для OS/2, MS-DOS и Windows, кроме того, можно было использовать компьютер, на котором они работали, в качестве рабочей станции. Неудачная рыночная судьба OS/2 не позволила системам LAN-Manager и LAN-Server захватить заметную долю рынка, но принципы работы этих сетевых систем во многом нашли свое воплощение в ОС 90-х годов – MS Windows NT.

В 80-е годы были приняты основные стандарты на коммуникационные технологии для локальных сетей: в 1980 г. – Ethernet, в 1985 г. – Token Ring, в конце 80-х – FDDI (Fiber Distributed Data Interface), распределенный интерфейс передачи данных по волоконно-оптическим каналам, двойное кольцо с маркером. Это позволило обеспечить совместимость сетевых ОС на нижних уровнях, а также стандартизировать операционные системы с драйверами сетевых адаптеров.

Для ПК применялись не только специально разработанные для них ОС (MS-Dos, NetWare, OS/2), но и адаптировались уже существующие ОС, в частности UNIX. Наиболее известной системой этого типа была версия UNIX компании Santa Cruz Operation (SCO UNIX).

В 90-е годы практически все операционные системы, занимающие заметное место на рынке, стали сетевыми. Сетевые функции встраиваются в ядро ОС, являясь ее неотъемлемой частью. В ОС используются средства мультиплексирования нескольких стеков протоколов, за счет которого компьютеры могут поддерживать одновременную работу с разнородными серверами и клиентами. Появились специализированные ОС, например, сетевая ОС IOS компании Cisco System, работающая в маршрутизаторах. Во второй половине 90-х годов все производители ОС усилили поддержку средств работы с интерфейсами. Кроме стека протоколов TCP/IP в комплект поставки начали включать утилиты, реализующие популярные сервисы Интернета: telnet, ftp, DNS, Web и др.

Особое внимание уделялось в последнем десятилетии и уделяется в настоящее время корпоративным сетевым операционным системам. Это одна из наиболее важных задач в обозримом будущем. Корпоративные ОС должны хорошо и устойчиво работать в крупных сетях, которые характерны для крупных организаций (предприятий, банков и т.п.), имеющих отделения во многих городах и, возможно, в разных странах. Корпоративная ОС должна без проблем взаимодействовать с ОС разного типа и работать на различных аппаратных платформах. Сейчас определились лидеры в классе корпоративных ОС – это MS Windows 2000/2003/XP Professional x64 Edition/7 Enterprise, UNIX и Linux-системы, а также Novell NetWare 6.5.

1.4. Назначение состав и функции ОС

В настоящее время существует большое количество различных типов операционных систем, отличающихся областями применения, аппаратными платформами, способами реализации и др. Назначение операционных систем можно разделить на четыре основные составляющие.

1. Организация (обеспечение) удобного интерфейса между приложениями и пользователями, с одной стороны, и аппаратурой компьютера – с другой. Вместо реальной

аппаратуры компьютера ОС представляет пользователю расширенную виртуальную машину, с которой удобнее работать и которую легче программировать. Вот список основных сервисов, предоставляемых типичными операционными системами.

1. *Разработка программ*: ОС представляет программисту разнообразные инструменты разработки приложений: редакторы, отладчики и т.п. Ему не обязательно знать, как функционируют различные электронные и электромеханические узлы и устройства компьютера. Часто пользователь не знает даже системы команд процессора, поскольку он может обойтись мощными высокоуровневыми функциями, которые предоставляет ОС.

2. *Исполнение программ*. Для запуска программы нужно выполнить ряд действий: загрузить в основную память программу и данные, инициализировать устройства ввода-вывода и файлы, подготовить другие ресурсы. ОС выполняет всю эту рутинную работу вместо пользователя.

3. *Доступ к устройствам ввода-вывода*. Для управления каждым устройством используется свой набор команд. ОС предоставляет пользователю единообразный интерфейс, который скрывает все эти детали и обеспечивает программисту доступ к устройствам ввода-вывода с помощью простых команд чтения и записи. Если бы программист работал непосредственно с аппаратурой компьютера, то для организации, например, чтения блока данных с диска ему пришлось бы использовать более десятка команд с указанием множества параметров. После завершения обмена программист должен был бы предусмотреть еще более сложный анализ результата выполненной операции.

4. *Контролируемый доступ к файлам*. При работе с файлами управление со стороны ОС предполагает не только глубокий учет природы устройства ввода-вывода, но и знание структур данных, записанных в файлах. Многопользовательские ОС, кроме того, обеспечивают механизм защиты при обращении к файлам.

5. *Системный доступ*. ОС управляет доступом к совместно используемой или общедоступной вычислительной системе в целом, а также к отдельным системным ресурсам. Она обеспечивает защиту ресурсов и данных от несанкционированного использования и разрешает конфликтные ситуации.

6. *Обнаружение ошибок и их обработка*. При работе компьютерной системы могут происходить разнообразные сбои за счет внутренних и внешних ошибок в аппаратном обеспечении, различного рода программных ошибок (переполнение, попытка обращения к ячейке памяти, доступ к которой запрещен и др.). В каждом случае ОС выполняет действия, минимизирующие влияние ошибки на работу приложения (от простого сообщения об ошибке до аварийной остановки программы).

7. *Учет использования ресурсов*. Хорошая ОС имеет средства учета использования различных ресурсов и отображения параметров производительности вычислительной системы. Эта информация важна для настройки (оптимизации) вычислительной системы с целью повышения ее производительности.

В результате реальная машина, способная выполнить только небольшой набор элементарных действий (машинных команд), с помощью операционной системы превращается в виртуальную машину, выполняющую широкий набор гораздо более мощных функций. Виртуальная машина тоже управляется командами, но уже командами более высокого уровня, например: удалить файл с определенным именем, запустить на выполнение прикладную программу, повысить приоритет задачи, вывести текст файла на печать и т.д. Таким образом, назначение ОС состоит в предоставлении пользователю (программисту) некоторой расширенной виртуальной машины, которую легче программировать и с которой легче работать, чем непосредственно с аппаратурой, составляющей реальный компьютер, систему или сеть.

2. Организация эффективного использования ресурсов компьютера. ОС не только представляет пользователям и программистам удобный интерфейс к аппаратным средствам компьютера, но и является своеобразным диспетчером ресурсов компьютера. К числу основных ресурсов современных вычислительных систем относятся процессоры, основная память, таймеры, наборы данных, диски, накопители на магнитных лентах (МЛ), внешние накопители памяти (CD/DVD/Blu-Ray/USB), принтеры, сетевые устройства и др. Эти ресурсы распределяются операционной системой между выполняемыми программами. В отличие от программы, которая

является статическим объектом, выполняемая программа – это динамический объект, он называется процессом и является базовым понятием современных ОС.

Управление ресурсами вычислительной системы с целью наиболее эффективного их использования является вторым назначением операционной системы. *Критерии эффективности*, в соответствии с которыми ОС организует управление ресурсами компьютера, могут быть различными. Например, в одних системах важен такой критерий, как пропускная способность вычислительной системы, в других – время ее реакции. Зачастую ОС должны удовлетворять нескольким, противоречащим друг другу критериям, что доставляет разработчикам серьезные трудности.

Управление ресурсами включает решение ряда общих, не зависящих от типа ресурса задач:

1. *планирование ресурса* – определение, какому процессу, когда и в каком качестве (если ресурс может выделяться частями) следует выделить данный ресурс;
2. *удовлетворение запросов на ресурсы* – выделение ресурса процессам;
3. *отслеживание состояния и учет использования ресурса* – поддержание оперативной информации о занятости ресурса и распределенной его доли;
4. *разрешение конфликтов между процессами*, претендующими на один и тот же ресурс.

Для решения этих общих задач управления ресурсами разные ОС используют различные алгоритмы, особенности которых, в конечном счете, определяют облик ОС в целом, включая характеристики производительности, область применения и даже пользовательский интерфейс. Таким образом, управление ресурсами составляют важное назначение ОС. В отличие от функций расширенной виртуальной машины большинство функций управления ресурсами выполняются операционной системой автоматически и прикладному программисту недоступны.

3. *Облегчение процессов эксплуатации аппаратных и программных средств вычислительной системы.* Ряд операционных систем имеет в своем составе наборы служебных программ, обеспечивающие резервное копирование, архивацию данных, проверку, очистку и дефрагментацию дисковых устройств и др.

Кроме того, современные ОС имеют достаточно большой набор средств и способов диагностики и восстановления работоспособности системы. Сюда относятся:

- диагностические программы для выявления ошибок в конфигурации ОС;
- средства восстановления последней работоспособной конфигурации;
- средства восстановления поврежденных и пропавших системных файлов и др.

Следует отметить еще одно назначение ОС.

4. *Возможность развития.* Современные ОС организуются таким образом, что допускают эффективную разработку, тестирование и внедрение новых системных функций, не прерывая процесса нормального функционирования вычислительной системы. Большинство операционных систем постоянно развиваются (нагляден пример Windows). Происходит это в силу следующих причин.

1. *Обновление и возникновение новых видов аппаратного обеспечения.* Например, ранние версии ОС UNIX и OS/2 не использовали механизмы страничной организации памяти (что это такое, мы рассмотрим позже), потому, что они работали на машинах, не обеспеченных соответствующими аппаратными средствами.

2. *Новые сервисы.* Для удовлетворения пользователей или нужд системных администраторов ОС должны постоянно предоставлять новые возможности. Например, может потребоваться добавить новые инструменты для контроля или оценки производительности, новые средства ввода-вывода данных (речевой ввод). Другой пример – поддержка новых приложений, использующих окна на экране дисплея.

3. *Исправления.* В каждой ОС есть ошибки. Время от времени они обнаруживаются и исправляются. Отсюда постоянные появления новых версий и редакций ОС. Необходимость регулярных изменений накладывает определенные требования на организацию операционных систем. Очевидно, что эти системы (как, впрочем, и другие сложные программы системы) должны иметь модульную структуру с четко определенными межмодульными связями (интерфейсами). Важную роль играет хорошая и полная документированность системы.

Перейдем к рассмотрению состава компонентов и функций ОС. Современные операционные системы содержат сотни и тысячи модулей (например, W2000 содержит 29 млн

строк исходного кода на языке С). Функции ОС обычно группируются либо в соответствии с типами локальных ресурсов, которыми управляет ОС, либо в соответствии со специфическими задачами, применимыми ко всем ресурсам. Совокупности модулей, выполняющих такие группы функций, образуют подсистемы операционной системы.

Наиболее важными подсистемами управления ресурсами являются подсистемы управления процессами, памятью, файлами и внешними устройствами, а подсистемами, общими для всех ресурсов, являются подсистемы пользовательского интерфейса, защиты данных и администрирования.

Управление процессами. Подсистема управления процессами непосредственно влияет на функционирование вычислительной системы. Для каждой выполняемой программы ОС организует один или более процессов. Каждый такой процесс представляется в ОС информационной структурой (таблицей, дескриптором, контекстом процессора), содержащей данные о потребностях процесса в ресурсах, а также о фактически выделенных ему ресурсах (область оперативной памяти, количество процессорного времени, файлы, устройства ввода-вывода и др.). Кроме того, в этой информационной структуре хранятся данные, характеризующие историю пребывания процесса в системе: текущее состояние (активное или заблокированное), приоритет, состояние регистров, программного счетчика и др.

В современных мультипрограммных ОС может существовать одновременно несколько процессов, порожденных по инициативе пользователей и их приложений, а также инициированных ОС для выполнения своих функций (системные процессы). Поскольку процессы могут одновременно претендовать на одни и те же ресурсы, подсистема управления процессами планирует очередность выполнения процессов, обеспечивает их необходимыми ресурсами, обеспечивает взаимодействие и синхронизацию процессов.

Управление памятью. Подсистема управления памятью производит распределение физической памяти между всеми существующими в системе процессами, загрузку и удаление программных кодов и данных процессов в отведенные им области памяти, настройку адресно-зависимых частей кодов процесса на физические адреса выделенной области, а также защиту областей памяти каждого процесса. Стратегия управления памятью складывается из стратегий выборки, размещения и замещения блока программы или данных в основной памяти. Соответственно используются различные алгоритмы, определяющие, когда загрузить очередной блок в память (по запросу или с упреждением), в какое место памяти его поместить и какой блок программы или данных удалить из основной памяти, чтобы освободить место для размещения новых блоков.

Одним из наиболее популярных способов управления памятью в современных ОС является виртуальная память. Реализация механизма виртуальной памяти позволяет программисту считать, что в его распоряжении имеется однородная оперативная память, объем которой ограничивается только возможностями адресации, предоставляемыми системой программирования.

Важная функция управления памятью – защита памяти. Нарушения защиты памяти связаны с обращениями процессов к участкам памяти, выделенной другим процессам прикладных программ или программ самой ОС. Средства защиты памяти должны пресекать такие попытки доступа путем аварийного завершения программы-нарушителя.

Управление файлами. Функции управления файлами сосредоточены в файловой системе ОС. Операционная система виртуализирует отдельный набор данных, хранящихся на внешнем накопителе, в виде файла – простой неструктурированной последовательности байтов, имеющих символическое имя. Для удобства работы с данными файлы группируются в каталоги, которые, в свою очередь, образуют группы – каталоги более высокого уровня. Файловая система преобразует символические имена файлов, с которыми работает пользователь или программист, в физические адреса данных на дисках, организует совместный доступ к файлам, защищает их от несанкционированного доступа.

Управление внешними устройствами. Функции управления внешними устройствами возлагаются на подсистему управления внешними устройствами, называемую также подсистемой ввода-вывода. Она является интерфейсом между ядром компьютера и всеми подключенными к нему устройствами. Спектр этих устройств очень обширен (принтеры, сканеры, мониторы, модемы, манипуляторы, сетевые адаптеры, АЦП разного рода и др.), сотни моделей этих

устройств отличаются набором и последовательностью команд, используемых для обмена информацией с процессором и другими деталями.

Программа, управляющая конкретной моделью внешнего устройства и учитывающая все его особенности, называется драйвером. Наличие большого количества подходящих драйверов во многом определяет успех ОС на рынке. Созданием драйверов занимаются как разработчики ОС, так и компании, выпускающие внешние устройства. ОС должна поддерживать четко определенный интерфейс между драйверами и остальными частями ОС. Тогда разработчики компаний-производителей устройств ввода-вывода могут поставлять вместе со своими устройствами драйверы для конкретной операционной системы.

Защита данных и администрирование. Безопасность данных вычислительной системы обеспечивается средствами отказоустойчивости ОС, направленными на защиту от сбоев и отказов аппаратуры и ошибок программного обеспечения, а также средствами защиты от несанкционированного доступа. Для каждого пользователя системы обязательна процедура логического входа, в процессе которой ОС убеждается, что в систему входит пользователь, разрешенный административной службой. Администратор вычислительной системы определяет и ограничивает возможности пользователей в выполнении тех или иных действий, т.е. определяет их права по обращению и использованию ресурсов системы.

Важным средством защиты являются функции аудита ОС, заключающегося в фиксации всех событий, от которых зависит безопасность системы. Поддержка отказоустойчивости вычислительной системы реализуется на основе резервирования (дисковые RAID-массивы, резервные принтеры и другие устройства, иногда резервирование центральных процессоров, в ранних ОС – дуальные и дуплексные системы, системы с мажоритарным органом и др.). Вообще обеспечение отказоустойчивости системы – одна из важнейших обязанностей системного администратора, который для этого использует ряд специальных средств и инструментов.

Интерфейс прикладного программирования. Прикладные программисты используют в своих приложениях обращения к операционной системе, когда для выполнения тех или иных действий им требуется особый статус, которым обладает только ОС. Возможности операционной системы доступны программисту в виде набора функций, который называется интерфейсом прикладного программирования (Application Programming Interface, API). Приложения обращаются к функциям API с помощью системных вызовов. Способ, которым приложение получает услуги операционной системы, очень похож на вызов подпрограмм.

Способ реализации системных вызовов зависит от структурной организации ОС, особенностей аппаратной платформы и языка программирования.

В ОС UNIX системные вызовы почти идентичны библиотечным процедурам. Ситуация в Windows иная (более подробно это рассмотрим далее).

Пользовательский интерфейс. ОС обеспечивает удобный интерфейс не только для прикладных программ, но и для пользователя (программиста, администратора). В ранних ОС интерфейс сводился к языку управления заданиями и не требовал терминала. Команды языка управления заданиями набивались на перфокарты, а результаты выполнения задания выводились на печатающее устройство.

Современные ОС поддерживают развитые функции пользовательского интерфейса для интерактивной работы за терминалами двух типов: алфавитно-цифрового и графического. При работе за алфавитно-цифровым терминалом пользователь имеет в своем распоряжении систему команд, развитость которой отражает функциональные возможности данной ОС. Обычно командный язык ОС позволяет запускать и останавливать приложения, выполнять различные операции с каталогами и файлами, получать информацию о состоянии ОС, администрировать систему. Команды могут вводиться не только в интерактивном режиме с терминала, но и считываться из так называемого командного файла, содержащего некоторую последовательность команд.

Программный модуль ОС, ответственный за чтение отдельных команд или же последовательности команд из командного файла, иногда называют командным интерпретатором (в MS-DOS – командным процессором).

Вычислительные системы, управляемые из командной строки, например UNIX-системы, имеют командный интерпретатор, называемый оболочкой (Shell). Она, собственно, не входит в состав ОС, но пользуется многими функциями операционной системы. Когда какой-либо

пользователь входит в систему, запускается оболочка. Стандартным терминалом для нее является монитор с клавиатурой. Оболочка начинает работу с печати приглашения (prompt) – знака доллара (или иного знака), говорящего пользователю, что оболочка ожидает ввода команды (аналогично управляется MS-DOS). Если теперь пользователь напечатает какую-либо команду, оболочка создает системный вызов и ОС выполнит эту команду. После завершения оболочка опять печатает приглашение и пытается прочесть следующую входную строку.

Ввод команд может быть упрощен, если операционная система поддерживает графический пользовательский интерфейс. В этом случае пользователь выбирает на экране нужный пункт меню или графический символ (так это происходит, например, в ОС Windows).

1.5. Архитектура операционной системы

Под архитектурой операционной системы понимают структурную и функциональную организацию ОС на основе некоторой совокупности программных модулей. В состав ОС входят исполняемые и объектные модули стандартных для данной ОС форматов, программные модули специального формата (например, загрузчик ОС, драйверы ввода-вывода), конфигурационные файлы, файлы документации, модули справочной системы и т.д.

На архитектуру ранних операционных систем обращалось мало внимания: во-первых, ни у кого не было опыта в разработке больших программных систем, а во-вторых, проблема взаимозависимости и взаимодействия модулей недооценивалась. В подобных монолитных ОС почти все процедуры могли вызывать одна другую. Такое отсутствие структуры было несовместимо с расширением операционных систем. Первая версия ОС OS/360 была создана коллективом из 5000 человек за 5 лет и содержала более 1 млн строк кода. Разработанная несколько позже операционная система Mastics содержала к 1975 году уже 20 млн строк. Стало ясно, что разработка таких систем должна вестись на основе модульного программирования.

Большинство современных ОС представляют собой хорошо структурированные модульные системы, способные к развитию, расширению и переносу на новые платформы. Какой-либо единой унифицированной архитектуры ОС не существует, но известны универсальные подходы к структурированию ОС. Принципиально важными универсальными подходами к разработке архитектуры ОС являются:

- модульная организация;
- функциональная избыточность;
- функциональная избирательность;
- параметрическая универсальность;
- концепция многоуровневой иерархической вычислительной системы, по которой ОС представляется многослойной структурой;
- разделение модулей на две группы по функциям: ядро – модули, выполняющие основные функции ОС, и модули, выполняющие вспомогательные функции ОС;
- разделение модулей ОС на две группы по размещению в памяти вычислительной системы: резидентные, постоянно находящиеся в оперативной памяти, и транзитные, загружаемые в оперативную память только на время выполнения своих функций;
- реализация двух режимов работы вычислительной системы: привилегированного режима (режима ядра – Kernel mode), или режима супервизора (supervisor mode), и пользовательского режима (user mode), или режима задачи (task mode);
- ограничение функций ядра (а следовательно, и количества модулей ядра) до минимального количества необходимых самых важных функций.

Первые ОС разрабатывались как монолитные системы без четко выраженной структуры (рис. 1.2).

Для построения монолитной системы необходимо скомпилировать все отдельные процедуры, а затем связать их вместе в единый объектный файл с помощью компоновщика (примерами могут служить ранние версии ядра UNIX или Novell NetWare). Каждая процедура видит любую другую процедуру (в отличие от структуры, содержащей модули, в которой большая часть информации является локальной для модуля, и процедуры модуля можно вызвать только через специально определенные точки входа).

Однако даже такие монолитные системы могут быть немного структурированными. При обращении к системным вызовам, поддерживаемым ОС, параметры помещаются в строго определенные места, такие как регистры или стек, а затем выполняется специальная команда прерывания, известная как вызов ядра или вызов супервизора. Эта команда переключает машину из режима пользователя в режим ядра, называемый также режимом супервизора, и передает управление ОС. Затем ОС проверяет параметры вызова, для того чтобы определить, какой системный вызов должен быть выполнен. После этого ОС индексирует таблицу, содержащую ссылки на процедуры, и вызывает соответствующую процедуру.

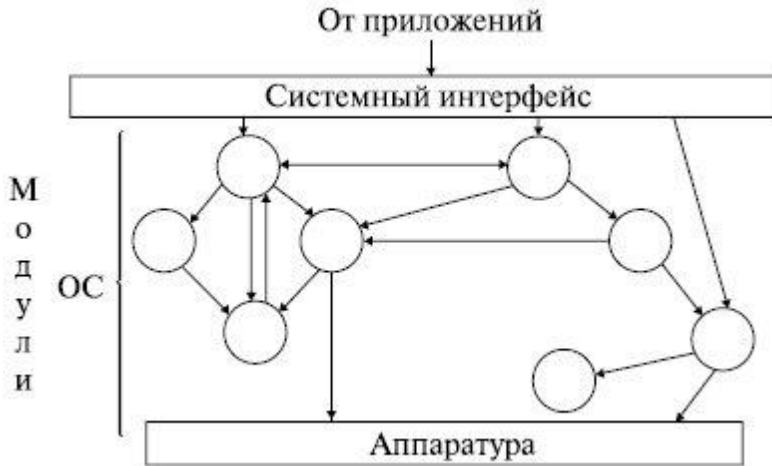


Рис. 1.2. Монолитная архитектура

Такая организация ОС предполагает следующую структуру:

- главная программа, которая вызывает требуемые сервисные процедуры;
- набор сервисных процедур, реализующих системные вызовы;
- набор утилит, обслуживающих сервисные процедуры.

В этой модели для каждого системного вызова имеется одна сервисная процедура. Утилиты выполняют функции, которые нужны нескольким сервисным процедурам. Это деление процедур на три слоя показано на рис. 1.3.

Классической считается архитектура ОС, основанная на концепции иерархической многоуровневой машины, привилегированном ядре и пользовательском режиме работы транзитных модулей. Модули ядра выполняют базовые функции ОС: управление процессами, памятью, устройствами ввода-вывода и т.п. Ядро составляет сердцевину ОС, без которой она является полностью неработоспособной и не может выполнить ни одну из своих функций. В ядре решаются внутрисистемные задачи организации вычислительного процесса, недоступные для приложения.

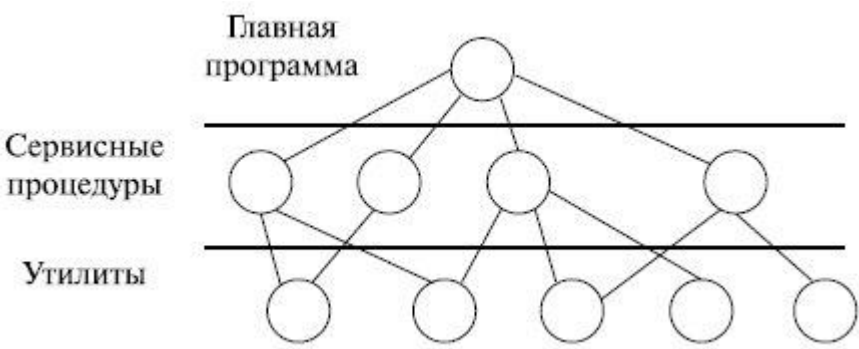


Рис. 1.3. Структурированная архитектура

Особый класс функций ядра служит для поддержки приложений, создавая для них так называемую прикладную программную среду. Приложения могут обращаться к ядру с запросами – системными вызовами – для выполнения тех или иных действий, например, открытие и чтение файла, получение системного времени, вывода информации на дисплей и т.д. Функции ядра, которые могут вызываться приложениями, образуют интерфейс прикладного программирования – API (Application Programming Interface).

Для обеспечения высокой скорости работы ОС модули ядра (по крайней мере, большая их часть) являются резидентными и работают в привилегированном режиме (Kernel mode). Этот режим, во-первых, должен обезопасить работу самой ОС от вмешательства приложений, и, во-вторых, должен обеспечить возможность работы модулей ядра с полным набором машинных инструкций, позволяющих собственно ядру выполнять управление ресурсами компьютера, в частности, переключение процессора с задачи на задачу, управлением устройствами ввода-вывода, распределением и защитой памяти и др.

Остальные модули ОС выполняют не столь важные функции, как ядро, и являются транзитными. Например, это могут быть программы архивирования данных, дефрагментации диска, сжатия дисков, очистки дисков и т.п.

Вспомогательные модули обычно подразделяются на группы:

- утилиты – программы, выполняющие отдельные задачи управления и сопровождения вычислительной системы;
- системные обрабатывающие программы – текстовые и графические редакторы (Paint, Imaging в Windows 2000), компиляторы и др.;
- программы предоставления пользователю дополнительных услуг (специальный вариант пользовательского интерфейса, калькулятор, игры, средства мультимедиа Windows 2000);
- библиотеки процедур различного назначения, упрощения разработки приложений, например, библиотека функций ввода-вывода, библиотека математических функций и т.п.

Эти модули ОС оформляются как обычные приложения, обращаются к функциям ядра посредством системных вызовов и выполняются в пользовательском режиме (user mode). В этом режиме запрещается выполнение некоторых команд, которые связаны с функциями ядра ОС (управление ресурсами, распределение и защита памяти и т.п.).

В концепции многоуровневой (многослойной) иерархической машины структура ОС также представляется рядом слоев. При такой организации каждый слой обслуживает вышележащий слой, выполняя для него некоторый набор функций, которые образуют межслойный интерфейс. На основе этих функций следующий верхний по иерархии слой строит свои функции – более сложные и более мощные и т.д. Такая организация системы существенно упрощает ее разработку, т.к. позволяет сначала "сверху вниз" определить функции слоев и межслойные интерфейсы, а при детальной реализации, двигаясь "снизу вверх", – наращивать мощность функции слоев. Кроме того, модули каждого слоя можно изменять без необходимости изменений в других слоях (но не меняя межслойных интерфейсов!).

Многослойная структура ядра ОС может быть представлена, например, вариантом, показанным на рис. 1.4.

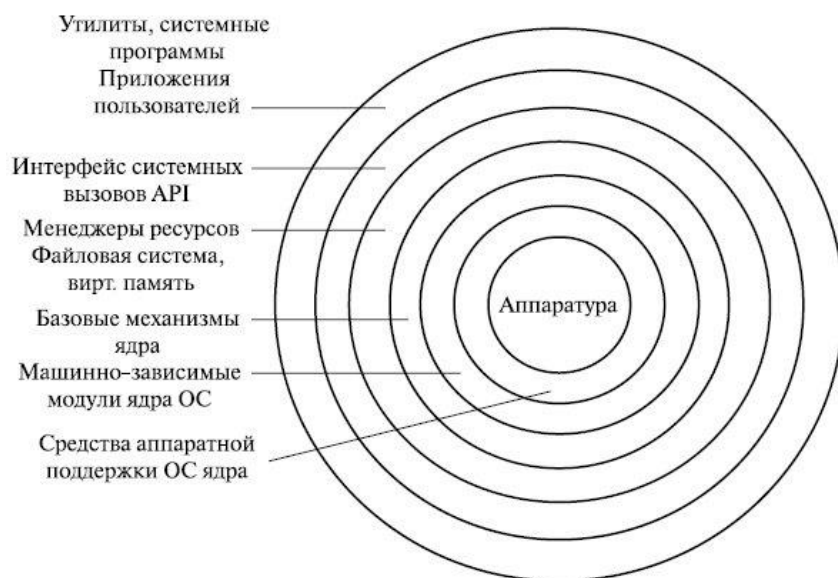


Рис. 1.4. Многослойная структура ОС

В данной схеме выделены следующие слои.

1. *Средства аппаратной поддержки ОС.* Значительная часть функций ОС может выполняться аппаратными средствами. Чисто программные ОС сейчас не существуют. Как правило, в современных системах всегда есть средства аппаратной поддержки ОС, которые прямо участвуют в организации вычислительных процессов. К ним относятся: система прерываний, средства поддержки привилегированного режима, средства поддержки виртуальной памяти, системный таймер, средства переключения контекстов процессов (информация о состоянии процесса в момент его приостановки), средства защиты памяти и др.

2. *Машинно-зависимые модули ОС.* Этот слой образует модули, в которых отражается специфика аппаратной платформы компьютера. Назначение этого слоя – "экранирование" вышележащих слоев ОС от особенностей аппаратуры (например, Windows 2000 – это слой HAL (Hardware Abstraction Layer), уровень аппаратных абстракций).

3. *Базовые механизмы ядра.* Этот слой модулей выполняет наиболее примитивные операции ядра: программное переключение контекстов процессов, диспетчерскую прерываний, перемещение страниц между основной памятью и диском и т.п. Модули этого слоя не принимают решений о распределении ресурсов, а только обрабатывают решения, принятые модулями вышележащих уровней. Поэтому их часто называют исполнительными механизмами для модулей верхних слоев ОС.

4. *Менеджеры ресурсов.* Модули этого слоя выполняют стратегические задачи по управлению ресурсами вычислительной системы. Это менеджеры (диспетчеры) процессов ввода-вывода, оперативной памяти и файловой системы. Каждый менеджер ведет учет свободных и используемых ресурсов и планирует их распределение в соответствии запросами приложений.

5. *Интерфейс системных вызовов.* Это верхний слой ядра ОС, взаимодействующий с приложениями и системными утилитами, он образует прикладной программный интерфейс ОС. Функции API, обслуживающие системные вызовы, предоставляют доступ к ресурсам системы в удобной компактной форме, без указания деталей их физического расположения.

Повышение устойчивости ОС обеспечивается переходом ядра в привилегированный режим. При этом происходит некоторое замедление выполнения системных вызовов. Системный вызов привилегированного ядра инициирует переключение процессора из пользовательского режима в привилегированный, а при возврате к приложению – обратное переключение. За счет этого возникает дополнительная задержка в обработке системного вызова (рис. 1.5). Однако такое решение стало классическим и используется во многих ОС (UNIX, VAX, VMS, IBM OS/390, OS/2 и др.).



Рис. 1.5. Обработка системного вызова

Многослойная классическая многоуровневая архитектура ОС не лишена своих проблем. Дело в том, что значительные изменения одного из уровней могут иметь трудно предвидимое влияние на смежные уровни. Кроме того, многочисленные взаимодействия между соседними уровнями усложняют обеспечение безопасности. Поэтому, как альтернатива классическому варианту архитектуры ОС, часто используется микроядерная архитектура ОС.

Суть этой архитектуры состоит в следующем. В привилегированном режиме остается работать только очень небольшая часть ОС, называемая микроядром. Микроядро защищено от остальных частей ОС и приложений. В его состав входят машинно-зависимые модули, а также модули, выполняющие базовые механизмы обычного ядра. Все остальные более высокоуровневые функции ядра оформляются как модули, работающие в пользовательском режиме. Так, менеджеры ресурсов, являющиеся неотъемлемой частью обычного ядра, становятся "периферийными" модулями, работающими в пользовательском режиме. Таким образом, в архитектуре с микроядром традиционное расположение уровней по вертикали заменяется горизонтальным. Это можно представить, как показано на рис. 1.6.

Внешние по отношению к микроядру компоненты ОС реализуются как обслуживающие процессы. Между собой они взаимодействуют как равноправные партнеры с помощью обмена сообщениями, которые передаются через микроядро. Поскольку назначением этих компонентов ОС является обслуживание запросов приложений пользователей, утилит и системных обрабатывающих программ, менеджеры ресурсов, вынесенные в пользовательский режим, называются серверами ОС, т.е. модулями, основным назначением которых является обслуживание запросов локальных приложений и других модулей ОС.



Рис. 1.6. Переход к микроядерной архитектуре

Схематично механизм обращений к функциям ОС, оформленным в виде серверов, выглядит, как показано на рис. 1.7.



Рис. 1.7. Клиент-серверная архитектура

Схема смены режимов при выполнении системного вызова в ОС с микроядерной архитектурой выглядит, как показано на рис. 1.8. Из рисунка ясно, что выполнение системного вызова сопровождается четырьмя переключениями режимов (4 t), в то время как в классической архитектуре – двумя. Следовательно, производительность ОС с микроядерной архитектурой при прочих равных условиях будет ниже, чем у ОС с классическим ядром.



Рис. 1.8. Обработка системного вызова в микроядерной архитектуре

В то же время признаны следующие достоинства микроядерной архитектуры [17]:

- единообразные интерфейсы;
- простота расширяемости;
- высокая гибкость;
- возможность переносимости;
- высокая надежность;
- поддержка распределенных систем;
- поддержка объектно-ориентированных ОС.

По многим источникам вопрос масштабов потери производительности в микроядерных ОС является спорным. Многое зависит от размеров и функциональных возможностей микроядра. Избирательное увеличение функциональности микроядра приводит к снижению количества переключений между режимами системы, а также переключений адресных пространств процессов.

Может быть, это покажется парадоксальным, но есть и такой подход к микроядерной ОС, как уменьшение микроядра.

Для возможности представления о размерах микроядер операционных систем в ряде источников приводятся такие данные:

- типичное микроядро первого поколения – 300 Кбайт кода и 140 интерфейсов системных вызовов;
- микроядро ОС L4 (второе поколение) – 12 Кбайт кода и 7 интерфейсов системных вызовов.

В современных операционных системах различают следующие виды ядер.

1. Наноядро (НЯ). Крайне упрощённое и минимальное ядро, выполняет лишь одну задачу – обработку аппаратных прерываний, генерируемых устройствами компьютера. После обработки посылает информацию о результатах обработки вышележащему программному обеспечению. НЯ используются для виртуализации аппаратного обеспечения реальных компьютеров или для реализации механизма гипервизора.

2. Микроядро (МЯ) предоставляет только элементарные функции управления процессами и минимальный набор абстракций для работы с оборудованием. Большая часть работы осуществляется с помощью специальных пользовательских процессов, называемых сервисами. В микроядерной операционной системе можно, не прерывая ее работы, загружать и выгружать новые драйверы, файловые системы и т. д. Микроядерными являются ядра ОС Minix и GNU Hurd и ядро систем семейства BSD. Классическим примером микроядерной

системы является Symbian OS. Это пример распространенной и отработанной микроядерной (а начиная с версии Symbian OS v8.1, и наноядерной) операционной системы.

3. Экзоядро (ЭЯ) предоставляет лишь набор сервисов для взаимодействия между приложениями, а также необходимый минимум функций, связанных с защитой: выделение и высвобождение ресурсов, контроль прав доступа и т. д. ЭЯ не занимается предоставлением абстракций для физических ресурсов – эти функции выносятся в библиотеку пользовательского уровня (так называемую libOS). В отличие от микроядра ОС, базирующиеся на ЭЯ, обеспечивают большую эффективность за счет отсутствия необходимости в переключении между процессами при каждом обращении к оборудованию.

4. Монолитное ядро (МНЯ) предоставляет широкий набор абстракций оборудования. Все части ядра работают в одном адресном пространстве. МНЯ требуют перекомпиляции при изменении состава оборудования. Компоненты операционной системы являются не самостоятельными модулями, а составными частями одной программы. МНЯ более производительны, чем микроядро, поскольку работает как один большой процесс. МНЯ является большинство Unix-систем и Linux. Монолитность ядер усложняет отладку, понимание кода ядра, добавление новых функций и возможностей, удаление ненужного, унаследованного от предыдущих версий кода. "Разбухание" кода монолитных ядер также повышает требования к объёму оперативной памяти.

5. Модульное ядро (Мод. Я) – современная, усовершенствованная модификация архитектуры МЯ. В отличие от "классических" МНЯ, модульные ядра не требуют полной перекомпиляции ядра при изменении состава аппаратного обеспечения компьютера. Вместо этого они предоставляют тот или иной механизм подгрузки модулей, поддерживающих то или иное аппаратное обеспечение (например, драйверов). Подгрузка модулей может быть как динамической, так и статической (при перезагрузке ОС после переконфигурирования системы). Мод. Я удобнее для разработки, чем традиционные монолитные ядра. Они предоставляют программный интерфейс (API) для связывания модулей с ядром, для обеспечения динамической подгрузки и выгрузки модулей. Не все части ядра могут быть сделаны модулями. Некоторые части ядра всегда обязаны присутствовать в оперативной памяти и должны быть жёстко "вшиты" в ядро.

6. Гибридное ядро (ГЯ) – модифицированные микроядра, позволяющие для ускорения работы запускать "несущественные" части в пространстве ядра. Имеют "гибридные" достоинства и недостатки. Примером смешанного подхода может служить возможность запуска операционной системы с монолитным ядром под управлением микроядра. Так устроены 4.4BSD и MkLinux, основанные на микроядре Mach. Микроядро обеспечивает управление виртуальной памятью и работу низкоуровневых драйверов. Все остальные функции, в том числе взаимодействие с прикладными программами, осуществляются монолитным ядром. Данный подход сформировался в результате попыток использовать преимущества микроядерной архитектуры, сохраняя по возможности хорошо отлаженный код монолитного ядра.

7. Наиболее тесно элементы микроядерной архитектуры и элементы монолитного ядра переплетены в ядре Windows NT. Хотя Windows NT часто называют микроядерной операционной системой, это не совсем так. Микроядро NT слишком велико (более 1 Мбайт), чтобы носить приставку "микро". Компоненты ядра Windows NT располагаются в вытесняемой памяти и взаимодействуют друг с другом путем передачи сообщений, как и положено в микроядерных операционных системах. В то же время все компоненты ядра работают в одном адресном пространстве и активно используют общие структуры данных, что свойственно операционным системам с монолитным ядром.

1.6. Классификация операционных систем

Все многообразие существующих (и ныне не использующихся) ОС можно классифицировать по множеству различных признаков. Остановимся на основных классификационных признаках.

1. По назначению ОС делятся на универсальные и специализированные. Специализированные ОС, как правило, работают с фиксированным набором программ (функциональных задач). Применение таких систем обусловлено невозможностью использования универсальной ОС по соображениям эффективности, надежности, защищенности и т.п., а также вследствие специфики решаемых задач.

Универсальные ОС рассчитаны на решение любых задач пользователей, но, как правило, форма эксплуатации вычислительной системы может предъявлять особые требования к ОС, т.е. к элементам ее специализации.

2. По способу загрузки можно выделить *загружаемые ОС* (большинство) и *системы, постоянно находящиеся в памяти* вычислительной системы. Последние, как правило, специализированные и используются для управления работой специализированных устройств (например, в БЦВМ баллистической ракеты или спутника, научных приборах, автоматических устройствах различного назначения и др.).

3. По особенностям *алгоритмов управления ресурсами*. Главным ресурсом системы является процессор, поэтому дадим классификацию по алгоритмам управления процессором, хотя можно, конечно, классифицировать ОС по алгоритмам управления памятью, устройствами ввода-вывода и т.д.

- Поддержка многозадачности (многопрограммности). По числу одновременно выполняемых задач ОС делятся на 2 класса: однопрограммные (однозадачные) – например, MS-DOS, MSX, и многопрограммные (многозадачные) – например, ОС ЕС ЭВМ, OS/360, OS/2, UNIX, Windows разных версий.

Однопрограммные ОС предоставляют пользователю виртуальную машину, делая более простым и удобным процесс взаимодействия пользователя с компьютером. Они также имеют средства управления файлами, периферийными устройствами и средства общения с пользователем. Многозадачные ОС, кроме того, управляют разделением совместно используемых ресурсов (процессор, память, файлы и т.д.), это позволяет значительно повысить эффективность вычислительной системы.

- Поддержка многопользовательского режима. По числу одновременно работающих пользователей ОС делятся: на однопользовательские (MS-DOS, Windows 3x, ранние версии OS/2) и многопользовательские (UNIX, Windows NT/2000/2003/XP/Vista).

Главное отличие многопользовательских систем от однопользовательских – наличие средств защиты информации каждого пользователя от несанкционированного доступа других пользователей. Следует заметить, что может быть однопользовательская мультипрограммная система.

- Виды многопрограммной работы. Специфику ОС во многом определяет способ распределения времени между несколькими одновременно существующими в системе процессами (или потоками). По этому признаку можно выделить 2 группы алгоритмов: не вытесняющая многопрограммность (Windows 3.x, NetWare) и вытесняющая многопрограммность (Windows 2000/2003/XP, OS/2, Unix).

В первом случае активный процесс выполняется до тех пор, пока он сам не отдает управление операционной системе. Во втором случае решение о переключении процессов принимает операционная система. Возможен и такой режим многопрограммности, когда ОС разделяет процессорное время между отдельными ветвями (потоками, волокнами) одного процесса.

- Многопроцессорная обработка. Важное свойство ОС – отсутствие или наличие средств поддержки многопроцессорной обработки. По этому признаку можно выделить ОС без поддержки мультипроцессорирования (Windows 3.x, Windows 95) и с поддержкой мультипроцессорирования (Solaris, OS/2, UNIX, Windows NT/2000/2003/XP).

Многопроцессорные ОС классифицируются по способу организации вычислительного процесса на асимметричные ОС (выполняются на одном процессоре, распределяя прикладные задачи по остальным процессорам) и симметричные ОС (децентрализованная система).

4. По области использования и форме эксплуатации. Обычно здесь выделяют три типа в соответствии с использованными при их разработке критериями эффективности:

- системы пакетной обработки (OS/360, ОС ЕС);
- системы разделения времени (UNIX, VMS);
- системы реального времени (QNX, RT/11).

Первые предназначались для решения задач в основном вычислительного характера, не требующих быстрого получения результатов. Критерий создания таких ОС – максимальная пропускная способность при хорошей загрузке всех ресурсов компьютера. В таких системах пользователь отстранен от компьютера.

Системы разделения времени обеспечивают удобство и эффективность работы пользователя, который имеет терминал и может вести диалог со своей программой.

Системы реального времени предназначены для управления техническими объектами (станок, спутник, технологический процесс, например доменный и т.п.), где существует предельное время на выполнение программ, управляющих объектом.

5. По аппаратной платформе (типу вычислительной техники), для которой они предназначаются, операционные системы делят на следующие группы.

- Операционные системы для смарт-карт. Некоторые из них могут управлять только одной операцией, например, электронным платежом. Некоторые смарт-карты являются JAVA-ориентированным и содержат интерпретатор виртуальной машины JAVA. Апплеты JAVA загружаются на карту и выполняются JVM-интерпретатором. Некоторые из таких карт могут одновременно управлять несколькими апплетами JAVA, что приводит к многозадачности и необходимости планирования.
- Встроенные операционные системы. Управляют карманными компьютерами (lialm OS, Windows CE – Consumer Electronics – бытовая техника), мобильными телефонами, телевизорами, микроволновыми печами и т.п.
- Операционные системы для персональных компьютеров, например, Windows 9.x, Windows XP, Linux, Mac OSX и др.
- Операционные системы мини-ЭВМ, например, RT-11 для PDP-11 – ОС реального времени, RSX-11 M для PDP-11 – ОС разделения времени, UNIX для PDP-7.
- Операционные системы мэйнфреймов (больших машин), например, OS/390, происходящая от OS/360 (IBM). Обычно ОС мэйнфреймов предполагает одновременно три вида обслуживания: пакетную обработку, обработку транзакций (например, работа с БД, бронирование авиабилетов, процесс работы в банках) и разделение времени.
- Серверные операционные системы, например, UNIX, Windows 2000, Linux. Область применения – ЛВС, региональные сети, Intranet, Internet.
- Кластерные операционные системы. Кластер – слабо связанная совокупность нескольких вычислительных систем, работающих совместно для выполнения общих приложений и представляющихся пользователю единой системой, например, Windows 2000 Cluster Server, Windows 2008 Server, Sun Cluster (базовая ОС – Solaris).

1.7. Эффективность и требования, предъявляемые к ОС

К операционным системам современных компьютеров предъявляется ряд требований. Главным требованием является выполнение основных функций эффективного управления ресурсами и обеспечения удобного интерфейса для пользователя и прикладных программ. Современная ОС должна поддерживать мультипрограммную обработку, виртуальную память, свопинг, развитый интерфейс пользователя (многооконный графический, аудио -, менюориентированный и т.д.), высокую степень защиты, удобство работы, а также выполнять многие другие необходимые функции и услуги. Кроме этих требований функциональной полноты, к ОС предъявляется ряд важных эксплуатационных требований.

1. **Эффективность.** Под эффективностью вообще любой технической (да и не только технической) системы понимается степень соответствия системы своему назначению, которая оценивается некоторым множеством показателей эффективности.

Поскольку ОС представляет собой сложную программную систему, она использует для собственных нужд значительную часть ресурсов компьютера. Часто эффективность ОС оценивают ее производительностью (пропускной способностью) – количеством задач

пользователей, выполняемых за некоторый промежуток времени, временем реакции на запрос пользователя и др.

На все эти показатели эффективности ОС влияет много различных факторов, среди которых основными являются архитектура ОС, многообразие ее функций, качество программного кода, аппаратная платформа (компьютер) и др.

2. Надежность и отказоустойчивость. Операционная система должна быть, по меньшей мере, так же надежна, как компьютер, на котором она работает. Система должна быть защищена как от внутренних, так и от внешних сбоев и отказов. В случае ошибки в программе или аппаратуре система должна обнаружить ошибку и попытаться исправить положение или, по крайней мере, постараться свести к минимуму ущерб, нанесенный этой ошибкой пользователям.

Надежность и отказоустойчивость ОС, прежде всего, определяются архитектурными решениями, положенными в ее основу, а также отлаженностью программного кода (основные отказы и сбои ОС в основном обусловлены программными ошибками в ее модулях). Кроме того, важно, чтобы компьютер имел резервные дисковые массивы, источники бесперебойного питания и др., а также программную поддержку этих средств.

3. Безопасность (защищенность). Ни один пользователь не хочет, чтобы другие пользователи ему мешали. ОС должна защищать пользователей и от воздействия чужих ошибок, и от попыток злонамеренного вмешательства (несанкционированного доступа). С этой целью в ОС как минимум должны быть средства аутентификации – определения легальности пользователей, авторизации – предоставления легальным пользователям установленных им прав доступа к ресурсам, и аудита – фиксации всех потенциально опасных для системы событий.

Свойства безопасности особенно важны для сетевых ОС. В таких ОС к задаче контроля доступа добавляется задача защиты данных, передаваемых по сети.

4. Предсказуемость. Требования, которые пользователь может предъявить к системе, в большинстве случаев непредсказуемы. В то же время пользователь предпочитает, чтобы обслуживание не очень сильно менялось в течение предположительного времени. В частности, запуская свою программу в системе, пользователь должен иметь основанное на опыте работы с этой программой приблизительное представление, когда ему ожидать выдачи результатов.

5. Расширяемость. В отличие от аппаратных средств компьютера полезная жизнь операционных систем измеряется десятками лет. Примером может служить ОС UNIX, да и MS-DOS. Операционные системы изменяются со временем, как правило, за счет приобретения новых свойств, например, поддержки новых типов внешних устройств или новых сетевых технологий. Если программный код модулей ОС написан таким образом, что дополнения и изменения могут вноситься без нарушения целостности системы, то такую ОС называют расширяемой. Операционная система может быть расширяемой, если при ее создании руководствовались принципами модульности, функциональной избыточности, функциональной избирательности и параметрической универсальности.

6. Переносимость. В идеальном случае код ОС должен легко переноситься с процессора одного типа на процессор другого типа и с аппаратной платформы (которые различаются не только типом процессора, но и способом организации всей аппаратуры компьютера) одного типа на аппаратную платформу другого типа. Переносимые ОС имеют несколько вариантов реализации для разных платформ, такое свойство ОС называется также многоплатформенностью. Достигается это свойство за счет того, что основная часть ОС пишется на языке высокого уровня (например С, С++ и др.) и может быть легко перенесена на другой компьютер (машинно-независимая часть), а некоторая меньшая часть ОС (программы ядра) является машинно-зависимой и разрабатывается на машинном языке другого компьютера.

7. Совместимость. Существует несколько "долгоживущих" популярных ОС (разновидности UNIX, MS-DOS, Windows 3.x, Windows NT, OS/2), для которых наработана широкая номенклатура приложений. Для пользователя, переходящего с одной ОС на другую, очень привлекательна возможность – выполнить свои приложения в новой операционной

системе. Если ОС имеет средства для выполнения прикладных программ, написанных для других операционных систем, то она совместима с этими системами. Следует различать совместимость на уровне двоичных кодов и совместимость на уровне исходных текстов. Кроме того, понятие совместимости включает также поддержку пользовательских интерфейсов других ОС.

8. *Удобство.* Средства ОС должны быть простыми и гибкими, а логика ее работы ясна пользователю. Современные ОС ориентированы на обеспечение пользователю максимально возможного удобства при работе с ними. Необходимым условием этого стало наличие у ОС графического пользовательского интерфейса и всевозможных мастеров – программ, автоматизирующих активизацию функций ОС, подключение периферийных устройств, установку, настройку и эксплуатацию самой ОС.

9. *Масштабируемость.* Если ОС позволяет управлять компьютером с различным числом процессоров, обеспечивая линейное (или почти такое) возрастание производительности при увеличении числа процессоров, то такая ОС является масштабируемой. В масштабируемой ОС реализуется симметричная многопроцессорная обработка. С масштабируемостью связано понятие кластеризации – объединения в систему двух (и более) многопроцессорных компьютеров. Правда, кластеризация направлена не столько на масштабируемость, сколько на обеспечение высокой готовности системы.

Следует заметить, что в зависимости от области применения конкретной операционной системы может изменяться и состав предъявляемых к ней требований.

Производители могут предлагать свои ОС в различных, различающихся ценой и производительностью, конфигурациях. Например, Microsoft продает [10]:

- Windows 2003 Server (до 4-х процессоров) – для малого и среднего бизнеса;
- Windows 2003 Advanced Server (до 8 процессоров, 2-узловой кластер) – для средних и крупных предприятий;
- Windows 2003 DataCenter Server (16-32 процессора, 4-узловой кластер) – для особо крупных предприятий.

1.8. Совместимость и множественные прикладные среды

В то время как многие архитектурные особенности ОС непосредственно касаются только системных программистов, концепция множественных прикладных (операционных) средств непосредственно связана с нуждами конечных пользователей – возможностью операционной системы выполнять приложения, написанные для других операционных систем. Такое свойство операционной системы называется совместимостью.

Совместимость приложений может быть на двоичном уровне и на уровне исходных текстов. Приложения обычно хранятся в ОС в виде исполняемых файлов, содержащих двоичные образы кодов и данных. Двоичная совместимость достигается в том случае, если можно взять исполняемую программу и запустить ее на выполнение в среде другой ОС.

Совместимость на уровне исходных текстов требует наличие соответствующего компилятора в составе программного обеспечения компьютера, на котором предполагается выполнить данное приложение, а также совместимости на уровне библиотек и системных вызовов. При этом необходима перекомпиляция исходных текстов приложения в новый исполняемый модуль.

Совместимость на уровне исходных текстов важна в основном для разработчиков приложений, в распоряжении которых эти исходные тексты имеются. Но для конечных пользователей практическое значение имеет только двоичная совместимость, так как только в этом случае они могут использовать один и тот же продукт в различных операционных системах и на различных машинах.

Вид возможной совместимости зависит от многих факторов. Самый главный из них – архитектура процессора. Если процессор применяет тот же набор команд (возможно, с добавлениями, как в случае IBM PC: стандартный набор + мультимедиа + графика + потоковые) и тот же диапазон адресов, то двоичная совместимость может быть достигнута достаточно просто. Для этого необходимо соблюдение следующих условий:

- API, который использует приложение, должен поддерживаться данной ОС;
- внутренняя структура исполняемого файла приложения должна соответствовать структуре исполняемых файлов данной ОС.

Если процессоры имеют разную архитектуру, то, кроме перечисленных условий, необходимо организовать эмуляцию двоичного кода. Например, широко используется эмуляция команд процессора Intel на процессоре Motorola 680x0 компьютера Macintosh. Программный эмулятор в этом случае последовательно выбирает двоичную инструкцию процессора Intel и выполняет эквивалентную подпрограмму, написанную в инструкциях процессора Motorola. Так как у процессора Motorola нет в точности таких же регистров, флагов, внутреннего АЛУ и др., как в процессорах Intel, он должен также имитировать (эмулировать) все эти элементы с использованием своих регистров или памяти.

Это простая, но очень медленная работа, поскольку одна команда Intel выполняется значительно быстрее, чем эмулирующая ее последовательность команд процессора Motorola. Выходом в таких случаях является применение так называемых прикладных программных сред или операционных сред. Одной из составляющих такой среды является набор функций интерфейса прикладного программирования API, который ОС предоставляет своим приложениям. Для сокращения времени на выполнение чужих программ прикладные среды имитируют обращение к библиотечным функциям.

Эффективность этого подхода связана с тем, что большинство сегодняшних программ работает под управлением GUI (графических интерфейсов пользователя) типа Windows, MAC или UNIX Motif, при этом приложения тратят 60-80% времени на выполнение функций GUI и других библиотечных вызовов ОС. Именно это свойство приложений позволяет прикладным средам компенсировать большие затраты времени, потраченные на покомандное эмулирование программ. Тщательно спроектированная программная прикладная среда имеет в своем составе библиотеки, имитирующие библиотеки GUI, но написанные на "родном" коде. Таким образом, достигается существенное ускорение выполнения программ с API другой операционной системы. Иначе такой подход называют трансляцией – для того, чтобы отличить его от более медленного процесса эмулирования по одной команде за раз.

Например, для Windows-программы, работающей на Macintosh, при интерпретации команд процессора Intel производительность может быть очень низкой. Но когда производится вызов функции GUI, открытие окна и др., модуль ОС, реализующий прикладную среду Windows, может перехватить этот вызов и перенаправить его на перекомпилированную для процессора Motorola 680x0 подпрограмму открытия окна. В результате на таких участках кода скорость работы программы может достичь (а, возможно, и превзойти) скорость работы на своём родном процессоре.

Чтобы программа, написанная для одной ОС, могла быть выполнена в рамках другой ОС, недостаточно лишь обеспечивать совместимость API. Концепции, положенные в основу разных ОС, могут входить в противоречия друг с другом. Например, в одной ОС приложению может быть разрешено управлять устройствами ввода-вывода, в другой – эти действия являются прерогативой ОС.

Каждая ОС имеет свои собственные механизмы защиты ресурсов, свои алгоритмы обработки ошибок и исключительных ситуаций, особую структуру процессора и схему управления памятью, свою семантику доступа к файлам и графический пользовательский интерфейс. Для обеспечения совместимости необходимо организовать бесконфликтное сосуществование в рамках одной ОС нескольких способов управления ресурсами компьютера.

Существуют различные варианты построения множественных прикладных сред, отличающиеся как особенностями архитектурных решений, так и функциональными возможностями, обеспечивающими разную степень переносимости приложений. Один из наиболее очевидных вариантов реализации множественных прикладных сред основывается на стандартной многоуровневой структуре ОС.

На рис. 1.9 ОС OS1 поддерживает кроме своих "родных" приложений приложения операционных систем OS2 и OS3. Для этого в её составе имеются специальные приложения, прикладные программные среды, которые транслируют интерфейсы "чужих" операционных систем API OS2 и API OS3 в интерфейс своей "родной" ОС – API OS1. Так, например, в случае

если бы в качестве OS2 выступала ОС UNIX, а в качестве OS1 – OS/2, для выполнения системного вызова создания процесса fork () в UNIX-приложении программная среда должна обращаться к ядру операционной системы OS/2 с системным вызовом DOS ExecPgm ().



Рис. 1.9. Организация множественных прикладных сред

К сожалению, поведение почти всех функций, составляющих API одной ОС, как правило, существенно отличается от поведения соответствующих функций другой ОС. Например, чтобы функция создания процесса в OS/2 Dos ExecPgm () полностью соответствовала функции создания процесса fork () в UNIX-подобных системах, её нужно было бы изменить и прописать новую функциональность: поддержку возможности копирования адресного пространства родительского процесса в пространство процесса-потомка.

Еще один способ построения множественных прикладных сред основан на микроядерном подходе. При этом очень важно отметить базовое, общее для всех прикладных сред отличие механизмов операционной системы от специфических для каждой из прикладных сред высокоуровневых функций, решающих стратегические задачи. В соответствии с микроядерной архитектурой все функции ОС реализуются микроядром и серверами пользовательского режима. Важно, что прикладная среда оформляется в виде отдельного сервера пользовательского режима и не включает базовых механизмов.

Приложения, используя API, обращаются с системными вызовами к соответствующей прикладной среде через микроядро. Прикладная среда обрабатывает запрос, выполняет его (возможно, обращаясь для этого за помощью к базовым функциям микроядра) и отправляет приложению результат. В ходе выполнения запроса прикладной среде приходится, в свою очередь, обращаться к базовым механизмам ОС, реализуемым микроядром и другими серверами ОС.

Такому подходу к конструированию множественных прикладных сред присущи все достоинства и недостатки микро ядерной архитектуры, в частности:

- очень просто можно добавлять и исключать прикладные среды, что является следствием хорошей расширяемости микро ядерных ОС;
- при отказе одной из прикладных сред остальные сохраняют работоспособность, что способствует надежности и стабильности системы в целом;
- низкая производительность микроядерных ОС сказывается на скорости работы прикладных средств, а значит, и на скорости работы приложений.

В итоге следует отметить, что создание в рамках одной ОС нескольких прикладных средств для выполнения приложений различных ОС представляет собой путь, который позволяет иметь единственную версию программы и переносить ее между различными операционными системами. Множественные прикладные среды обеспечивают совместимость на двоичном уровне данной ОС с приложениями, написанными для других ОС.

1.9. Виртуальные машины как современный подход к реализации множественных прикладных сред

Понятие "монитор виртуальных машин" (МВМ) возникло в конце 60-х годов как программный уровень абстракции, разделявший аппаратную платформу на несколько виртуальных машин. Каждая из этих виртуальных машин (ВМ) была настолько похожа на базовую физическую машину, что существующее программное обеспечение могло выполняться на ней в неизменном виде. В то время вычислительные задачи общего характера решались на дорогих мэйнфреймах (типа IBM/360), и пользователи высоко оценили способность МВМ распределять дефицитные ресурсы среди нескольких приложений.

В 80-90-е годы существенно снизилась стоимость компьютерного оборудования и появились эффективные многозадачные ОС, что уменьшило ценность МВМ в глазах пользователей. Мэйнфреймы уступили место мини-компьютерам, а затем ПК, и нужда в МВМ отпала. В результате из компьютерной архитектуры попросту исчезли аппаратные средства для их эффективной реализации. К концу 80-х в науке и на производстве МВМ воспринимались не иначе как исторический курьез [10].

Сегодня МВМ – снова в центре внимания. Корпорации Intel, AMD, Sun Microsystems и IBM создают стратегии виртуализации, в научных лабораториях и университетах для решения проблем мобильности, обеспечения безопасности и управляемости развиваются подходы, основанные на виртуальных машинах. Что же произошло между отставкой МВМ и их возрождением?

В 90-е годы исследователи из Стэнфордского университета начали изучать возможность применения ВМ для преодоления ограничений оборудования и операционных систем. Проблемы возникли у компьютеров с массовой параллельной обработкой (Massively Parallel Processing, MPP), которые плохо поддавались программированию и не могли выполнять имеющиеся ОС. Исследователи обнаружили, что с помощью виртуальных машин можно сделать эту неудобную архитектуру достаточно похожей на существующие платформы, чтобы использовать преимущества готовых ОС. Из этого проекта вышли люди и идеи, ставшие золотым фондом компании VMware (www.vmware.com), первого поставщика МВМ для компьютеров массового применения.

Как ни странно, развитие современных ОС и снижение стоимости оборудования привели к появлению проблем, которые исследователи надеялись решить с помощью МВМ. Дешевизна оборудования способствовала быстрому распространению компьютеров, но они часто бывали недогруженными, требовали дополнительных площадей и усилий по обслуживанию. А следствиями роста функциональных возможностей ОС стали их неустойчивость и уязвимость.

Чтобы уменьшить влияние системных аварий и защититься от взломов, системные администраторы вновь обратились к однозадачной вычислительной модели (с одним приложением на одной машине). Это привело к дополнительным расходам, вызванным повышенными требованиями к оборудованию. Перенос приложений с разных физических машин на ВМ и консолидация этих ВМ на немногих физических платформах позволили повысить эффективность использования оборудования, снизить затраты на управление и производственные площади. Таким образом, способность МВМ к мультиплексированию аппаратных средств – на этот раз во имя консолидации серверов и организации коммунальных вычислений – снова возродила их к жизни.

В настоящее время МВМ стал не столько средством организации многозадачности, каким он был когда-то задуман, сколько решением проблем обеспечения безопасности, мобильности и надежности. Во многих отношениях МВМ дает создателям операционных систем возможность развития функциональности, невозможной в нынешних сложных ОС. Такие функции, как миграция и защита, намного удобнее реализовать на уровне МВМ, поддерживающих обратную совместимость при развертывании инновационных решений в области операционных систем при сохранении предыдущих достижений.

Виртуализация – развивающаяся технология. В общих словах, виртуализация позволяет отделить ПО от нижележащей аппаратной инфраструктуры. Фактически она разрывает связь между определенным набором программ и конкретным компьютером. Монитор виртуальных машин отделяет программное обеспечение от оборудования и формирует промежуточный уровень

между ПО, выполняемым виртуальными машинами, и аппаратными средствами. Этот уровень позволяет МВМ полностью контролировать использование аппаратных ресурсов *гостевыми операционными системами (GuestOS)*, которые выполняются на ВМ.

МВМ создает унифицированное представление базовых аппаратных средств, благодаря чему физические машины различных поставщиков с разными подсистемами ввода-вывода выглядят одинаково и ВМ выполняются на любом доступном оборудовании. Не заботясь об отдельных машинах с их тесными взаимосвязями между аппаратными средствами и программным обеспечением, администраторы могут рассматривать оборудование просто как пул ресурсов для оказания любых услуг по требованию.

Благодаря *полной инкапсуляции* состояния ПО на ВМ монитор МВМ может отобразить ВМ на любые доступные аппаратные ресурсы и даже перенести с одной физической машины на другую. Задача балансировки нагрузки в группе машин становится тривиальной, и появляются надежные способы борьбы с отказами оборудования и наращивания системы. Если нужно отключить отказавший компьютер или ввести в строй новый, МВМ способен соответствующим образом перераспределить виртуальные машины. Виртуальную машину легко тиражировать, что позволяет администраторам по мере необходимости оперативно предоставлять новые услуги.

Инкапсуляция также означает, что администратор может в любой момент приостановить или возобновить работу ВМ, а также сохранить текущее состояние виртуальной машины либо вернуть ее к предыдущему состоянию. Располагая возможностью универсальной отмены, удастся легко справиться с авариями и ошибками конфигурации. Инкапсуляция является основой обобщенной модели мобильности, поскольку приостановленную ВМ можно копировать по сети, сохранять и транспортировать на сменных носителях.

МВМ играет роль посредника во всех взаимодействиях между ВМ и базовым оборудованием, поддерживая выполнение множества виртуальных машин на единой аппаратной платформе и обеспечивая их надежную изоляцию. МВМ позволяет собрать группу ВМ с низкими потребностями в ресурсах на отдельном компьютере, снизив затраты на аппаратные средства и потребность в производственных площадях.

Полная изоляция также важна для надежности и обеспечения безопасности. Приложения, которые раньше выполнялись на одной машине, теперь можно распределить по разным ВМ. Если одно из них в результате ошибки вызовет аварию ОС, другие приложения будут от нее изолированы и продолжат работу. Если же одному из приложений угрожает внешнее нападение, атака будет локализована в пределах "скомпрометированной" ВМ. Таким образом, МВМ – это инструмент реструктуризации системы для повышения ее устойчивости и безопасности, не требующий дополнительных площадей и усилий по администрированию, которые необходимы при выполнении приложений на отдельных физических машинах.

МВМ должен связать аппаратный интерфейс с ВМ, сохранив полный контроль над базовой машиной и процедурами взаимодействия с ее аппаратными средствами. Для достижения этой цели существуют разные методы, основанные на определенных технических компромиссах. При поиске таких компромиссов принимаются во внимание основные требования к МВМ: совместимость, производительность и простота. Совместимость важна потому, что главное достоинство МВМ – способность выполнять унаследованные приложения. Производительность определяет величину накладных расходов на виртуализацию – программы на ВМ должны выполняться с той же скоростью, что и на реальной машине. Простота необходима, поскольку отказ МВМ приведет к отказу всех ВМ, выполняющихся на компьютере. В частности, для надежной изоляции требуется, чтобы МВМ был свободен от ошибок, которые злоумышленники могут использовать для разрушения системы.

Вместо того чтобы заниматься сложной переработкой кода гостевой операционной системы, можно внести некоторые изменения в основную операционную систему, изменив некоторые наиболее "мешающие" части ядра. Подобный подход называется паравиртуализацией [10]. Ясно, что в этом случае адаптировать ядро ОС может только автор, и, например, Microsoft не проявляет желания адаптировать популярное ядро Windows 2000 к реалиям конкретных виртуальных машин.

При паравиртуализации разработчик МВМ переопределяет интерфейс виртуальной машины, заменяя непригодное для виртуализации подмножество исходной системы команд более

удобными и эффективными эквивалентами. Заметим, что хотя ОС нужно портировать для выполнения на таких ВМ, большинство обычных приложений могут выполняться в неизменном виде.

Самый большой недостаток паравиртуализации – несовместимость. Любая операционная система, предназначенная для выполнения под управлением паравиртуализованного монитора МВМ, должна быть портирована в эту архитектуру, для чего нужно договариваться о сотрудничестве с поставщиками ОС. Кроме того, нельзя использовать унаследованные операционные системы, а существующие машины не удастся легко заменить виртуальными.

Чтобы добиться высокой производительности и совместимости при виртуализации архитектуры x86, компания VMware разработала новый метод виртуализации, который объединяет традиционное прямое выполнение с быстрой трансляцией двоичного кода "на лету". В большинстве современных ОС режимы работы процессора при выполнении обычных прикладных программ легко поддаются виртуализации, а следовательно, их можно виртуализировать посредством прямого выполнения. Непригодные для виртуализации привилегированные режимы может выполнять транслятор двоичного кода, исправляя "неудобные" команды x86. В результате получается высокопроизводительная виртуальная машина, которая полностью соответствует оборудованию и поддерживает полную совместимость ПО.

Преобразованный код очень похож на результаты паравиртуализации. Обычные команды выполняются в неизменном виде, а команды, нуждающиеся в специальной обработке (такие как ROPF и команды чтения регистров сегмента кода), транслятор заменяет последовательностями команд, которые подобны требующимся для выполнения на паравиртуализованной виртуальной машине. Однако есть важное различие: вместо того, чтобы изменять исходный код операционной системы или приложений, транслятор двоичного кода изменяет код при его выполнении в первый раз.

Хотя трансляция двоичного кода требует некоторых дополнительных расходов, при нормальных рабочих нагрузках они незначительны. Транслятор обрабатывает лишь часть кода, и скорость выполнения программ становится сопоставимой со скоростью прямого выполнения – как только заполнится кэш-память трассировки.

Трансляция двоичного кода также помогает оптимизировать прямое выполнение. Например, если при прямом выполнении привилегированного кода часто происходит перехват команд, это может привести к существенным дополнительным расходам, поскольку при каждом перехвате управление передается от виртуальной машины к монитору и обратно. Трансляция кода может устранить многие из таких перехватов, что приведет к снижению накладных расходов на виртуализацию. Это особенно верно для центральных процессоров с длинными конвейерами команд, в частности, для современного семейства x86, в котором перехват связан с высокими дополнительными расходами.

1.10. Эффекты виртуализации

Экспертиза современных продуктов и недавние исследования раскрывают некоторые интересные возможности развития МВМ и требования, которые они предъявляют к технологиям виртуализации.

Администраторы центра данных могут с единой консоли быстро вводить в действие ВМ и управлять тысячами виртуальных машин, выполняющихся на сотнях физических серверов. Вместо того чтобы конфигурировать отдельные компьютеры, администраторы будут создавать по имеющимся шаблонам новые экземпляры виртуальных серверов и отображать их на физические ресурсы в соответствии с политиками администрирования. Уйдет в прошлое взгляд на компьютер как на средство предоставления конкретных услуг. Администраторы будут рассматривать компьютеры просто как часть пула универсальных аппаратных ресурсов (примером тому может служить виртуальный центр VMware VirtualCenter).

Отображение виртуальных машин на аппаратные ресурсы очень динамично. Возможности миграции работающих ВМ (подобные тем, которые обеспечивает технология VMotion компании VMware) позволяют ВМ быстро перемещаться между физическими машинами в соответствии с потребностями центра данных. МВМ сможет справляться с такими традиционными проблемами,

как отказ оборудования, за счет простого перемещения ВМ с отказавшего компьютера на исправный. Возможность перемещения работающих ВМ облегчит решение аппаратных проблем, таких как планирование профилактического обслуживания, окончание срока действия лизингового договора и модернизация оборудования: администраторы станут устранять эти проблемы без перерывов в работе.

Еще недавно нормой являлась ручная миграция, но сейчас уже распространены инфраструктуры виртуальных машин, которые автоматически выполняют балансировку нагрузки, прогнозируют отказы аппаратных средств и соответствующим образом перемещают ВМ, создают их и уничтожают в соответствии со спросом на конкретные услуги.

Решение проблем на уровне МВМ положительно сказывается на всех программах, выполняющихся на ВМ, независимо от их возраста (унаследованная или новейшая) и поставщиков. Независимость от ОС избавляет от необходимости покупать и обслуживать избыточную инфраструктуру. Например, из нескольких версий ПО службы поддержки или резервного копирования останется лишь одна – та, которая работает на уровне МВМ.

Виртуальные машины сильно изменили отношение к компьютерам. Уже сейчас простые пользователи умеют легко создавать, копировать и совместно использовать ВМ. Модели их применения значительно отличаются от привычных, сложившихся в условиях вычислительной среды с ограниченной доступностью аппаратных средств. А разработчики ПО могут применять такие продукты, как VMware Workstation, чтобы легко установить компьютерную сеть для тестирования или создать собственный набор испытательных машин для каждой цели.

Повышенная мобильность ВМ значительно изменила способы их применения. Такие проекты, как Collective и Internet Suspend/Resume, демонстрируют возможность перемещения всей вычислительной среды пользователя по локальной и территориально-распределенной сети. Доступность высокоемких недорогих сменных носителей, например, жестких дисков USB, означает, что потребитель может захватить свою вычислительную среду с собой, куда бы он ни направлялся.

Динамический характер компьютерной среды на базе ВМ требует и более динамичной топологии сети. Виртуальные коммутаторы, виртуальные брандмауэры и оверлейные сети становятся неотъемлемой частью будущего, в котором логическая вычислительная среда отделится от своего физического местоположения.

Виртуализация обеспечивает высокий уровень работоспособности и безопасности благодаря нескольким ключевым возможностям.

Локализация неисправностей. Большинство отказов приложений происходят из-за ошибок ПО. Виртуализация обеспечивает логическое разделение виртуальных разделов, поэтому программный сбой в одном разделе никак не влияет на работу приложения в другом разделе. Логическое разделение также позволяет защищаться от внешних атак, что повышает безопасность консолидированных сред.

Гибкая обработка отказов. Виртуальные разделы можно настроить так, чтобы обеспечить автоматическую обработку отказов для одного или нескольких приложений. Благодаря средствам обеспечения высокой степени работоспособности, заложенным сейчас в платформы на базе процессоров Intel® Itanium® 2 и Intel® Xeon™ MP, требуемый уровень услуг часто можно обеспечить, предусмотрев аварийный раздел на той же платформе, где работает основное приложение. Если требуется еще более высокий уровень работоспособности, аварийный раздел можно разместить на отдельной платформе.

Разные уровни безопасности. Для каждой виртуальной машины можно установить разные настройки безопасности. Это позволит ИТ-организациям обеспечить высокий уровень контроля за конечными пользователями, а также гибкое распределение административных привилегий.

МВМ имеют мощный потенциал для реструктуризации существующих программных систем в целях повышения уровня защиты, а также облегчают развитие новых подходов к построению безопасных систем. Сегодняшние ОС не обеспечивают надежной изоляции, оставляя машину почти незащищенной. Перемещение механизмов защиты за пределы ВМ (чтобы они выполнялись параллельно с ОС, но были изолированы от нее) позволяет сохранить их функциональные возможности и повысить устойчивость к нападениям.

Размещение средств безопасности за пределами ВМ – привлекательный способ изоляции сети. Доступ к сети предоставляется ВМ после проверки, гарантирующей, что она, с одной стороны, не представляет угрозы, а с другой – неуязвима для нападения. Управление доступом к сети на уровне ВМ превращает виртуальную машину в мощный инструмент борьбы с распространением злонамеренного кода.

Мониторы МВМ особенно интересны в плане управления многочисленными группами программ с различными уровнями безопасности. Благодаря отделению ПО от оборудования ВМ обеспечивают максимальную гибкость при поиске компромисса между производительностью, обратной совместимостью и степенью защиты. Изоляция программного комплекса в целом упрощает его защиту. В сегодняшних ОС почти невозможно судить о безопасности отдельного приложения, поскольку процессы плохо изолированы от друг друга. Таким образом, безопасность приложения зависит от безопасности всех остальных приложений на машине.

Гибкость управления ресурсами, которую обеспечивают МВМ, может сделать системы более стойкими к нападениям. Возможность быстро тиражировать ВМ и динамически адаптироваться к большим рабочим нагрузкам станет основой мощного инструмента, позволяющего справиться с нарастающими перегрузками из-за внезапного наплыва посетителей на Web-сайте или атаки типа "отказ в обслуживании".

Модель распространения программных продуктов на основе ВМ потребует от поставщиков ПО корректировки лицензионных соглашений. Лицензии на эксплуатацию на конкретном процессоре или физической машине не приживутся в новых условиях, в отличие от лицензий на число пользователей или неограниченных корпоративных лицензий. Пользователи и системные администраторы будут отдавать предпочтение операционным средам, которые легко и без особых затрат распространяются в виде виртуальных машин.

Возрождение МВМ существенно изменило представления разработчиков программных и аппаратных средств о структурировании сложных компьютерных систем и управлении ими. Кроме того, МВМ обеспечивают обратную совместимость при развертывании инновационных решений в области операционных систем, которые позволяют решать современные задачи, сохраняя предыдущие достижения. Эта их способность станет ключевой при решении грядущих компьютерных проблем.

Виртуализация предоставляет также преимущества для сред разработки и тестирования ПО. Различные этапы цикла создания ПО, включая получение рабочей версии, можно выполнять в разных виртуальных разделах одной и той же платформы. Это поможет повысить степень полезного использования аппаратного обеспечения и упростить управление жизненным циклом. Во многих случаях ИТ-организации получают возможность тестировать новые и модернизированные решения на имеющихся рабочих платформах, не прерывая производственный процесс. Это не только упрощает миграцию, но также позволяет сократить расходы, устранив необходимость дублирования вычислительной среды.

Освобождая разработчиков и пользователей от ресурсных ограничений и недостатков интерфейса, виртуальные машины снижают уязвимость системы, повышают мобильность программного обеспечения и эксплуатационную гибкость аппаратной платформы.

Компьютерные системы существуют и продолжают развиваться благодаря тому, что разработаны по законам иерархии и имеют хорошо определенные интерфейсы, отделяющие друг от друга уровни абстракции. Использование таких интерфейсов облегчает независимую разработку аппаратных и программных подсистем силами разных групп специалистов. Абстракции скрывают детали реализации нижнего уровня, уменьшая сложность процесса проектирования.

Подсистемы и компоненты, разработанные по спецификациям разных интерфейсов, не способны взаимодействовать друг с другом. Например, приложения, распространяемые в двоичных кодах, привязаны к определенной ISA и зависят от конкретного интерфейса к операционной системе. Несовместимость интерфейсов может стать сдерживающим фактором, особенно в мире компьютерных сетей, в котором свободное перемещение программ столь же необходимо, как и перемещение данных.

Виртуализация позволяет обойти эту несовместимость. Виртуализация системы или компонента (например, процессора, памяти или устройства ввода/вывода) на конкретном уровне

абстракции отображает его интерфейс и видимые ресурсы на интерфейс и ресурсы реальной системы. Следовательно, реальная система выступает в роли другой, виртуальной системы или даже нескольких виртуальных систем.

В отличие от абстракции, виртуализация не всегда нацелена на упрощение или сокрытие деталей. Например, при отображении виртуальных дисков на реальные программные средства виртуализации используют абстракцию файла как промежуточный шаг. Операция записи на виртуальный диск преобразуется в операцию записи в файл (и следовательно, в операцию записи на реальный диск). Отметим, что в данном случае никакого абстрагирования не происходит – уровень детализации интерфейса виртуального диска (адресация секторов и дорожек) ничем не отличается от уровня детализации реального диска.