

ССЫЛКИ НА ВСЕ ИЗВЕСТНЫЕ МНЕ СТАРЫЕ(И НЕ ОЧЕНЬ) ДОКИ НА ВТОРОЙ СТРАНИЦЕ			СТАВЬ ПОСОСЬ ЕСЛИ ЗВЕЗДОЧКА	Самая важная док
если что-то не ознакомили				Обязательное рассмотрение
Как раскрыты квадратик:	ОК, редактирование ячейки закрыто	потоковый скачок	<- рекомендовано к ознакомлению	Скачать - займал!
	необходимо доделать	конспекты у2018 не очень		FAQ по материалам второй курс
	не просмотрено	гуглопапка		дока препода 34-35 (есть вопросы, но много багов)
В ТАБЛИЧКЕ ТЕПЕРЬ РАЗРЕШЕНО ТОЛЬКО КОММЕНТИРОВАНИЕ				
ОЧЕНЬ ИНТЕРЕСНЫЕ ВОПРОСЫ:				
Рассказать про альтернативные подходы к увеличению скорости доступа, кроме кэша (нет, это не треды, и не НВМ) когда нам выгодно чтобы халязы случались чаще?				Какое-то название файла(1)
				Какое-то название файла(2)
Почему ячейки квадратные? (0009)				конспекты Желе
Потому что такimoto( наиболее подходящий ответ: хотим уменьшить длину проводов -> хотим уменьшить периметр ячейки -> хотим оптимизировать площадь ->>> делаем ячейку квадратной) Ибо потому что это упрощает производство? Оптимально расположить на плате квадратик ИМХО проще чем кругиком, звездочка, кук. Ну этого до, но не на приговорили? Скажем Сказал, что он из, почему они именно квадратные, они вполне себе могут быть и прямоугольными, всё зависит от жаланий производителей				Скачать 2u17
0.0 География				Билеты by u2017
Что лучше: Уругвай или Парагвай?			Уругвай	еще билеты by u2017
1.1 Элементная база вычислительной системы: логические элементы, триггеры				
1 Нарисовать логичный сумматор на заданном элементе	схема с Mux, где S - бит суммы по модулю, C - бит переноса			
2 Почему у триггер не может быть асинхронным?	Потому что если убрать сигнал синхронизации, то в состоянии 1 1 триггер будет постоянно менять значение на выходах			
3 Нарисовать RS, D, JK, T триггеры. Их синхронные версии, если есть	Вот тут есть описание всех триггеров			
4 Нарисовать логичный сумматор	схема, где S - бит суммы по модулю, C - бит переноса			
5 Что есть (диз)мультиплексор. Нарисовать заданной ёмкости	Статья простого человека, которую приятно читать			
6 Нарисовать побитовый счетчик	http://www.fhio.ru/wiki/index.php?title=Bitwise Counter.jpg			
7 Рисовать логические элементы на полках (мб бинофильных) транзисторах	да			
8 Как выглядит RS? Что будет, если хотим сделать RS на AND, а не на OR	ну придется еще выходные сигналы инвертировать, и всё ок, выглядит так же, как RS на NOR'ах, только вместо NOR - NAND			
9 RS-триггер на И-НЕ	Безопасно Тестиров			
10 Триггеры: Зачем нужна синхронизация?	Так как электрический импульс доходит от выходов одного элемента до входов другого не мгновенно, возможна ситуация, когда мы подали на логический элемент две единицы, одна из них дошла раньше. Несинхронизированный логический элемент вернет результат от неверных входных данных. Это грустно, поэтому нужна синхронизация. Прием такт синхронизации должен быть достаточно длинным, чтобы все значения успели установиться на входах элементов			
11 Построить умножитель (на входе A0, A1, A2, B0, B1, B2, которые задают 2 трибитах числа, на выходе Q0-Q5 - результат их произведения)	Матричный умножитель			
12 Чем отличается двухтранзисторный инвертор от одотранзисторного?	В двухтранзисторном инвертере выход подключен либо к земле, либо к питанию. В одотранзисторном инвертере при подаче высокого напряжения выход будет не к чему не подключен. Это очень плохо, поскольку в таком случае можем получить какое-либо случайное напряжение от других транзисторов. Высокое энергопотребление и тепловыделение. Ну, надо вставлять резистор, чтобы не было короткого замыкания. Такая конструкция постоянно жрет энергию	да		
13 Нарисуй JK. Объясни как сменить переключением инвертора	схематический пример с описанием			Здесь есть текст самого величайшего и он веселый, стр 14-15, лучше вставить
1.2 Оперативная память: статическая/динамическая, организация.				
1 Почему ячейки хранятся матрично, а не в линию.	Потому что много проводов - глюки. Хранение в линии 1Мб уже дает нам больше 10^6 проводов. А если хранить матрично, то всего около 10^3. Выбор очевиден			
2 Двухпортовая ячейка стат. и динам. памяти. Схемы ячеек памяти	ссылка			
3 Что идеальное среднее между однопортовой и двухпортовой ячейками статической памяти?	Многобанковая память.			
4 Преимущества, недостатки динамической и статической ячеек памяти.	Динамическая: - дешевле - занимает меньше места, чем статическая (транзистор + конденсатор вместо 8 или 6 транзисторов) - необходимо постоянно перезаписывать из-за утечки заряда из конденсаторов - необходимо перезаписывать после чтения, т.к информация при чтении теряется Статическая: - не нужно перезаписывать - работает быстрее - занимает много места - дорогая в производстве			Регенерация динамической памяти: - программист (давно дано) - контроллер памяти (он выполняет команды процессора и физическую память) - модуль памяти (лучший вариант, он не забывает шину своим глюком)
5 Как лучше опрашивать память? Почему номер строки и столбца лучше передавать последовательно, а не параллельно?	Потому что в два раза меньше проводов. Без потери скорости. Так как память устроена матрично нет смысла передавать одновременно потому что нужно время для открытия строки. Номер столбца будет использоваться не сразу. Поэтому можно передавать последовательно уменьшая при этом шину.			
6 Когда бывает выгодно отправлять адрес строки и столбца одновременно?	Если у нас SRAM, то нам не нужно записывать строку в буфер, а значит можно одновременно передавать строку и столбец для ускорения			
7 Что общего между ячейкой однопортовой и двухпортовой памяти?	они в целом практически одинаково устроены. Только немножко дополнительных транзисторов для второго набора проводов			
8 Может ли динамическая память быть двухпортовой?	да (см. 2 вопрос)			
9 А для чего используется правый провод (CSOL)?	Почему существует линия пол CSOL? Предположим, что мы убрали транзистор справа и линию пол CSOL. Тогда, например, если мы хотим записать значение 1 в ячейку, а в ней у нас значение 0, то ячейка инвертор будет мешать нам сделать это, ведь он будет пытаться вернуть значение 0. Если же у нас будет присутствовать линия пол CSOL, мы можем легче подавить старое значение на инверторе, ведь минимам состояние сразу обоих инверторов. Если мы хотим прогнать значение из ячейки без линии пол CSOL, нам придется "честно" смотреть на значение на линии COL, а сигнал может быть недостаточно силен, чтобы отпугить 0 от 1. А обычно мы смотрим на разницу в напряжении между COL и пол COL.			
10 Нарисовать иррегулярный модуль памяти с триггерами на 8 бит	Такое есть упрощенно в плане схемы			
11 Нарисуй ячейку SDRAM и DRAM. Что за транзисторы в них?	да			
12 Зачем записывать строку перед открытием колоды?	Чтобы не потерять данные (вспомним, что при чтении конденсаторы разряжаются)			
13 Почему не козем особо двухпортовую память?	Потому что есть многобанковая, а она дешевле и проще. Возможная параллельность не равносильна удорожанию стоимости и увеличению размеров ячеек.			Note: Двухпортовую козем когда надо уметь параллельно читать и писать. Например видеопластини(Cisco RAM). Одновременно считывают что показывать на мониторе и меняют чтоб картинка двигалась
1.3 Оперативная память: характеристики, типы динамической памяти, NUMA.				
1 Есть две оперативки, 8-8-8-20 и 9-8-8-20, какая лучше?	При случайных обращениях к памяти - одинаковы. Если читаем несколько столбцов из одной строки, то первая лучше. Случайные обращения к памяти: IPR + IRAS - два последних числа В одной строке: CL - первое число			
2 Дробный тайминг	Вспомним, что DDR работает 2 раза за такт, то есть возможна ситуация, что мы начали на фронте, а закончили на спаде, получается, что прошло какое-то целое число тактов + 0.5 такта.			В DDR мы увеличили внутреннюю шину и теперь у нас первая половина данных выплывается сразу, а вторая ждет такта. Это же и есть дробный тайминг?
3 Что общего в DDR? Почему это лучше, чем увеличить частоту в два раза? что изменили в DDR2 и DDR3?	Передаем информацию два раза за такт. Лучше, потому что увеличение частоты требует больше энергии. В ddr2 увеличили частоту внешней шины и ширину внутренней, уменьшили напряжение питания. В ddr3 сделали все это еще раз.			
4 При каких улучшениях RAM улучшалась скорость передачи? Зачем записывать строку?	1) IPR, EDO, BEDO, SDRAM, DDR, DDR2, DDR3 2) Чем чище строки их содержание попадает во внутренний буфер, при этом склеивается само содержимое строки. При закрытии строки информация записывается из буфера в саму строку. (см ответ 1.2.12)			
5 Как мы получаем в ddr два раза инфу за такт?	Увеличили ширину внутренней шины. Теперь достаем из памяти в 2 раза больше информации, кладем во внутренний буфер, передаем первую половину, затем вторую.			
6 Почему не получится просто ставить вездь выше скорость передачи, а не только в видеокартах?	В чем нем в видеокартах. Им важна скорость передачи данных, а скорость доступа не так важна. Видеокарточки могут позволить себе отводить большие тепла. Забивая на тепловыделение мы можем поднять скорость передачи(подробней в конспекте почтате)			
7 Есть ли какая-то организация памяти многопроцессорной не NUMA?	известно, UMA (Uniform Memory Access)			
8 Нарисовать многопроцессорную не NUMA архитектуру	он давал на картинке			
9 От чего зависит скорость передачи и скорость доступа?	Скорость доступа это время между тем когда мы подали команду и получением реакции на эту команду. Равно IRCD+HCL. Зачастую мы просим не один кусочек информации, а несколько подряд. У памяти мы просим прочитать не один байт, а дать кусок памяти с такого то места и длиной столько то. Память передает нам один кусочек и через время передает второй, третий и т.д. Скорость передачи данных это время через которое мы получим второй кусок информации после получения первого. Например винчестер. Ему нужно время чтобы найти информацию, которую мы запросили. Это скорость доступа. Теперь когда он нашел данные и выдает порядк кусок данных, которые мы запросили, ему не надо не задерживать. Получается что скорость передачи данных существенно быстрее скорости доступа.			
10 Почему синхронизация в SDRAM - хорошо?	Мем в том, что раньше тактовые частоты на памяти и на остальных устройствах не совпадали, поэтому приходилось ждать столько тактов, сколько нужно и еще чуть-чуть, чтобы гарантированно получить корректные данные. Потом умные чуваки решили посадить все устройства на одну и ту же тактовую частоту, поэтому теперь можно ждать ровно столько, сколько нужно			
11 DDR, мем он лучше, почему просто не увеличить частоту у SDRAM?	Если сравнивать только скорость передачи, то немем, но надо понимать, что если увеличить частоту в 2 раза, это приведет к увеличению энергопотребления. То есть в DDR не происходит увеличение энергопотребления, а в SDRAM с удвоенной частотой произойдет бы			
12 Изменения в DDR2, DDR3, DDR4?	DDR2: Увеличили ширину внутренней шины до 256 бит и частоту внешней шины => удвоение скорости передачи. Уменьшили напряжение питания до 1.8В DDR3: То же самое, удвоили ширину внутренней до 512 бит и частоту внешней => опять удвоение скорости передачи. Уменьшили напряжение питания до 1.5В DDR4: Удвоили кол-во банков памяти (до 16) + улучшили энергоэффективность (опять уменьшили напряжение питания, -1.3V), чуть поднови частоты, но не засчет удвоения внутренней шины			Скажем не рассказывать про удвоение банков, и точно говорить что скорость передачи не увеличилась
13 Отличие SDRAM, BEDO и DDR	И Точно такой вопрос? Может сразу все выложить, начиная с FPM? так			
14 Что такое NUMA? Нарисуй.	NUMA - NonUniformMemoryAccess. Устройство памяти компьютера в мультипроцессорной системе. При таком устройстве скорость доступа к памяти зависит от ее расположения относительно процессора. Контроллер памяти входит в процессор (важно) Пример NUMA-архитектуры на картинке (обор). У каждого процессора быстрый доступ к своей памяти, более долгий к памяти соседних процессоров и еще более долгий к памяти процессора расположенного по диагонали			заставля
15 Основные характеристики оперативки				
CL - IRCD - IPR - IRAS CL - CAS (Column Address Strobe) Latency - число тактов между отправкой адреса столбца и началом получения данных. IRCD - RAS (Row Address Strobe) to CAS Delay - минимальное число тактов, которое нужно подождать при открытой строке, прежде чем отчитать столбец. CL + IRCD - время доступа при закрытой строке IPR - Row Precharge - минимальное число тактов, которое нужно подождать между закрытием строки и открытием новой. CL + IRCD + IPR - время доступа при открытой неправильной строке IRAS - Row Active Strobe - минимальное число тактов между открытием и закрытием строки.				
16 Что за буфер в EDO памяти?	буфер адреса столбца. Быстрее считываем столбец, если они из одной строки			
17 Чем GDDR отличается от DDR? Как мы их можем сравнить?	GDDR оптимизирована на максимальную скорость передачи данных, но она дорогая, много жрет энергии, хуже скорость доступа			
18 Какие характеристики RAM влияют на скорость работы кэша?	скорость передачи + ширина шины + куда в DDR стали выкидывать данные два раза за такт, кэш стало приятно (кэш-линия в 64Кб стала заполняться за 4 такта, а не за 4)			Какой то посох. Кашу на самом деле до лампы на скорость передачи(вроде). А вот скорость доступа ему сильно нужнее. Это не точно но с этим надо разбираться
1.4 Кэш-память.				
1 Уровни кэша	Обычно есть три уровня кэша: L1, L2 и L3. L3 обычно общий для всех ядер. L1 обычно делится на L1i (кэш для инструкций) и L1d (кэш для данных)			

2	Когда каш бывает вреден	Такая ситуация возможна, если у нас каш с write-allocate (т.е. при каш-промежке при записи мы сначала подгружаем каш-линию в каш, а потом уже пишем туда значение). В таком случае, если у нас случайные обращения к памяти, то мы будем делать две операции с памятью (чтение и запись) на каждый запрос. А без каш или с no-write-allocate (при промажке при записи - просто пишем в память), мы бы делали только одну операцию (запись), это интеллектуальный ответ, ответ по SKYW выглядит так: "если мы используем данные, лежащие в разных линиях (например с шагом в размер линии), тогда мы передаем все каш линии (пишем 64 байт)"		
3	Что такое ассоциативность. Чем плохи и хороши 0 и 1 ассоциативности.	см. вопрос 7		
4	Когда мы записываем в память значение из каш при Write-Back	Когда выжидаем каш-линию.		
5	Бывают ли алгоритмы, для которых безразлично - есть каш или нет?	Да, те, которые читают данные, расположенные в случайных местах памяти Пример: сортировка кучей	Алгоритму не может быть безразлично наличие каш: он либо полезен либо вреден, как он может не влиять на скорость? Ну и да, пример алгоритмов плес	"любой алгоритм с более-менее случайным доступом к памяти" (и) сканов
6	Зачем нужен каш; типичные размеры кашей разных уровней	Каждый раз чтение из памяти долго (чтение занимает порядка 200 тактов). Поэтому давайте заведем небольшую быструю память поближе к процессору. 1-ый уровень - 2x32 Кб и -4 такта, 2-ой уровень - (256-512) Кб - 12 тактов, 3-ий уровень - 6-8 Мб и -10 тактов	Почему доступ к памяти именно 200 тактов занимает, это можно как-то высчитать? Это стоило спрашивать на паре, инфр про 200 тактов бы Скаков и он дронил на нас 100	
7	Что такое ассоциативность? Рассказать про типичные ассоциативности, адресацию в них, преимущества и недостатки	Ассоциативность каш - это когда мы разделяем каш на п. множества, в каждом из которых хранится л. линий. Блок в памяти сначала соотносится с множеством, а потом записывается в любую каш-линию этого множества. 1. ассоциативность (direct mapping) - в каждом множестве ровно одна каш-линия. У такого подхода к организации каш есть преимущества: быстрая ищба данных, высокая аппаратная часть (сравнение только одних байт). И недостатки: низкий hit gauge. Т.е. если мы часто достаем из памяти данные, которые попадают в одно множество, то нам грустно: каш-линия у нас всего одна и данные будут постоянно заменять друг друга. Выигрыш в скорости собою не поглотит. Полная ассоциативность: если можем писать любой блок памяти в любую каш-линию. Очень долгий и дорогой поиск по кашу. В итоге как-то пытаемся балансировать между двумя крайностями. Обычно делают 8-ассоциативность.	Знаете ли вы, что обычно 8-ассоциативность. А что там по адресации? <- в конспектах почитать, там написано почему картинку, чтобы объяснить	
8	Как процессор воспринимает каш? Видит он его или нет?	Контролер памяти в процессоре общается с памятью через шину памяти и сам по себе не знает, что именно отвечает на его запросы. На этой шине могут сидеть другие устройства: другие процессоры, контроллеры со стороны оперативной памяти и, в том числе, каш-подсистема. На уровне каш зашифрованные абсолютно прозрачны. Однако с каш-подсистемой можно общаться, если она это поддерживает. В современных системах детально можно либо попросить каш-подсистему что-нибудь зашифровать, однако это не значит, что она моментально исполнит этот запрос (это рекомендация, а каш-подсистема может быть свои приоритеты). Более того, она вообще не обязана слушать эти рекомендации и может не исполнить рекомендации ищбы. Выбавить, откуда взять или куда сохранить данные явно - нехорошо. С точки зрения "горящего работы" для процессора есть заметные различия: обмен данными при наличии каш в разы быстрее.		
9	Какой недостаток look-aside?	Если нужны данные найдены в каш, то нужно посылать сигнал отмены в память. В современном мире каш-подсистема случается примерно в 90 случаях. А значит, жерная каша между кашей и CPU дала бы больше выигрыша в скорости, чем немого более быстрый доступ к памяти. А в look-aside сделать жерную шину не получится.	White-back: Второй минус Wb: в разе нужно записывать в память, а думал каш дублирует память, разве нет? Явно нет, там write-back, он часто хранит "грязные" данные (т.е. не совпадающие с данными в памяти)	ну и зачем ты совсем стер то, что было?
10	Сравни write-through и write-back.	Write-through: - процесс реализовывать - забиваем кашу запросами к памяти	и пишем в память только кашу вытесняем линию из каш - сканов реализовывать (т.е. нужно хранить каш-линия содержит измененные данные памяти или нет) - если делаем запрос на чтение и промываемся, то зачастую обращаемся к памяти дважды: сначала чтобы записать вытесненную строку в память, а затем чтобы нужные данные (все равно лучше, чем write-through)	
11	Что быстрее write-through или write-back?	WB, так как кашу у нас в память притягивает запрос на чтение не зашифрованных данных, перед тем не будет микротрафик очереди из бессмысленных дублирующих запросов на запись		
12	В чем плюсы и минусы инкапсуляции и экспозиции архитектуры?	Если все каш-линии каш присутствуют в кашах более высокого уровня - то этот каш является инкапсулированным. Если каш-линии каш не присутствуют в кашах более высокого уровня, то он экспозиционный.	Экспозиционная (AMD) Каш-линия хранится только в одном каше + быстрее добавляем в каш + тратим меньше места + дольше ищба	Инкапсулиная (Intel) Каш-линия хранится в более чем одном каше + быстрее выжидаем из каш (не нужно проговаривать в кашах большего уровня, т.е. там уже присутствуют каша этой каш-линии) + быстрее ищба + тратим больше места
<b>1.5 Протоколы когерентности каш-памяти</b>				
1	Зачем нужны протоколы когерентности? Как обходятся без них?	Представим себе ситуацию, что у нас несколько ядер на процессоре. У каждого каша есть свои иеш L1, L2 и общий иеш L3 (примерно так всегда на современных процессорах), думаю, что одно ядро знает данные из RAM и изменяет их, другое ядро тоже хочет вить данные по тому же адресу и может взять уже старую, неактуальную версию этих данных. По той причине и существует протоколы когерентности, которые контролируют такие баги. Обойтись без этих протоколов нельзя, кроме как сделать один иеш на все процессоры (см.88 строку)	Можно обойтись без когерентности еще одним способом. Когда у нас идет взаимодействие между двумя разными транзакциями, то давайте будем обращаться тем нужным данным в каждом транзе, когда они там используются. Еще можно вывести иеш за шину, но тогда прирываем по скорости инфр из прилегающего ядра	5
2	Как заработать состояние Owned (исключительный каш)?	MOESI: Пусть к каше подключены два каша со своими кашей. Первое каш отправляет на шину памяти запрос на чтение иешей X. Так как этот иешей нет ни в одном из кашей (ока в состоянии Invalid), то запрос отправляется в оперативную память. Когда оттуда приходит ответ, в каше первого каша иешей X появляется в состоянии Exclusive (а каше второго каша она Invalid, поэтому зашифрована в единственном экземпляре). Затем первое каш должно отправить на шину памяти запрос на запись в иешей X. Эта иешей есть в его каше, поэтому он обрабатывает данный запрос: записывает в ней новые данные и переводит каш в состояние Modified. Так как она присутствует только в каше первого каша, то в остальных кашах ничего не изменяется. Данное второе каш отправляет на шину памяти запрос на чтение иешей X. В его каше эта иешей находится в состоянии Invalid, однако в каше первого каша она присутствует в состоянии Modified, поэтому каш первого каша отправляет на шину иеш значение, а сам переводит ее в себе в состояние Owned. После этого каш второго каш принимает ответ от каш первого, отдает его второму кашу, и записывает иешей X себе в состоянии Shared.	(* Где-то тут нужно указание на картинку с переходами между состояниями в MOESI, которая на листе в карандашам *)	
3	Рассказать про состояние "E" в MESI	E - Exclusive. Данный блок присутствует только в текущем каше и является "чистым". Может быть изменен на Shared при получении запроса на чтение. Может быть изменен на Modified при получении запроса на запись в этот блок.		
4	Процессор хочет записать в память строку MESI (все случаи) в MOESI (все пометы в O)	MESI: 1. Попади в Modified - изменением локально и все ок 2. Exclusive - переводим в состояние Modified, пишем локально, все ок 3. Shared - промыв все каш вытеснив свои каши, переводим в состояние Modified, пишем, все ок 4. Invalid - промыв все каш вытеснив свои каши, читаем себе в каш, изменяем новое состояние строки - Modified	MOESI: Очистив - пишем данные локально, переводим в состояние Modified, сообщаем остальным кашам удалить свои копии этой каш-линии	
5	В чем отличие пай от пайе, чем пайе пай	MSI блок тем, что мы забиваем шину ненужными запросами. А именно: каждый раз при изменении состояния с Shared на Modified мы промыв другие каш вытеснив данные (ведут они у них есть), но чаще всего данные есть только у нас в каше. Давайте избавимся от этой проблемы, добавим состояние Exclusive. Теперь мы можем изменить с E на M без каких-либо оповещений других кашей.		
6	Какие ошибки могут возникать в системе на двух каш при отсутствии когерентности? Прямые каш, кашы без когерентности	ири есть два транза: Первый транза переключает а = 5 на кашу и читается while a = 5 Второй транза и и читается while a = 5 С одной стороны, а не может быть одновременно равна и не равна 5, а с другой, без когерентности именя та и получается. Если добавим когерентность, то закончим только второй.		
7	Какие каш исправлять не используют когерентность? Как эти каш исправляет когерентность?	При когерентности первый транз получит значение а = 6 и цикл прервется. Не используем когерентность - в тех местах программы, где происходит обмен данными между транзакциями хотя бы потенциально может происходить адрес записи коммунукаций "обратного иеш"		
8	В каких случаях когерентность не возникает при транза нехот?	Если данные, которыми пользуются каши, не пересекаются (Тут надо сказать, что разные КЕШ-ЛИНИИ, т.е. я могу одним иешом попросить записать в адрес 228, а другим в 229, то я двумя иешами попал в одну каш-линию и проиграл, котк, по факту, адреса разные)		
9	Можно ли придумать такое устройство КЕШ, для которого не нужны протоколы когерентности	Да, один общий иеш на все процессоры		
10	Почему проблему когерентности не решает WT кашей?	усты у нас есть транза и 2 каш соответственно, а оба из них подгружено а = 3, тогда если один транз изменит значение а, скажем а = 5, то да, если каш WT другой каш не заметит это изменение пока не выскочит свой записи а не подгрузит его заново. То есть даже при WT в каше будет лежать не актуальное значение		
<b>1.6 Носители информации: магнитные, оптические и на основе флеш-памяти, RAID</b>				
1	SLC и MLC, что такое, + и -	SLC - Single Level Cell - есть два уровня напряжения: высокие и низкие. Храним один бит в клетке. MLC - Multi Level Cell - есть четыре уровня напряжения. Храним два бита в одной клетке. SLC: + выше скорость записи + большая надежность/долговечность + ниже энергопотребление + дороже MLC: + выше плотность записи за счт того, что храним 2 бита в одном транзисторе + стоимость ниже + ниже скорость записи и чтения		
2	Когда не работает TRIM	из прошлого вопроса: 1) если это RAID 2) если подключен по USB (ограничение протокола) 3) если диск зашифрован (тогда нужны соседние блоки для дешифрования) 4) не поддерживается ОС 5) не поддерживают устройство 6) если диск переполнен	еще про TRIM	у TRIM есть много проблем с RAID. Для исправления этого, ядро выпускает специальный пайе, но он только для иешей.
3	Сравнение скорости и SSD. Почему скорости могут плавать у обоих, и почему?	У винчестера разное кол-во секторов в начале(краю) и в конце(центре), в угловая скорость - сопна, поэтому скорость вначале помет будет даже в 2 раза больше чем в конце, а у SSD если мы записываем на чистый сектор, то стирать не надо, а если перезаписываем(стираем + записываем), то будет много дольше. Версия 2: Сравнение скорости: скорость доступа у SSD выше. Скорость передачи у HDD выше. Сам Паша говорит, что скорость чтения(записи)(пишем поверх без речистки) у HDD выше. Плавать скорость у HDD может например: 1) при чтении(записи наверх/тоже) при увеличении номера сектора(при движении от края к центру) скорость чтения снижается. 2) можем перезаписать не по 4 Кбайт, а меньше. Тогда должны прочитать весь сектор, оставить ту-то-то-то фрагмент, переписать остаток CRC/ECC и опять записать 3) плохо разбиваем HDD на кластеры/нада, чтобы раздел на диске начинался с сектора, номер которого, кратен 8 (4Кбайт(физическое разбиение на сектора) / 512 байт(логический размер сектора, эквивалент промывке 512б) / размер кластера у файловой системы кратен 4 Кбайт) у SSD: 1) зависит от режима транзистора/SLC MLC TLC GLC в порядке падения скорости записи, чтения) 2) скорость записи плавают от "чистоты" SSD (на новый все пишется быстро) 3) крутость контроллера 4) кол-во свободных иешей для маневров контроллера(если время переставит, чистим(ищем больше места, там быстрее работает SSD, т.е. контроллер может быстрее обработать запроса(на числе повышения скорости доступа), когда(ищем) завершит свои деления)	что значит плавать? Не имеет постоянного времени доступа? именно	
4	Достоинства и недостатки флеш-памяти. Почему время поиска у flash меньше, чем у жестких дисков	Потому что производители ищбы и сделали их туло параллельными, т.е. де-факто у нас не один, а много таких дисков		
5	Что такое trim?	Чел сверху(снизу) прав, но параллелизм это про скорость передачи, а поиск == скорость доступа, она у flash меньше, потому что там не надо "попадать в сектор" то есть покусить диск + передвинуть голову, там просто есть адрес памяти по которому мы идём		
6	Что такое write amplification и почему оно может быть больше/меньше 1. Стандартные значения WA	TRIM - команда, которая позволяет ОС сообщить SSD какие блоки данных уже не используются и их можно удалить. Появилась скоро после появления SSD, т.е. скорость записи в ненужный блок у SSD очень высокая, (надо сначала стереть, потом записать)		Хорошо записано
7	Почему размер сектора жесткого диска равно 4096?	Вспомним, что в SSD есть интеллектуальный контроллер, который равномерно распределяет нагрузку, то есть если он видит, что у какого-то иешей часто записываем, а в каку-то нет, то очередной запрос пойдет в ту, в которую рекомендациями он пойдет каша должна, но физически в ту область которая менее нагружена, потому что контроллер перемещает), чтобы не потерять данные которые были в редко используемой иешей, мы должны их тоже перемещать, откуда и появляется + 1 Версия 2: Есть хороший перекресток: -> через иешейку! Если кратко, поверх записывать не можем. Но минимальный блок для стирания гораздо больше минимального блока для записи. Хотим поменять небольшой блок файла. Тогда: должны прочитать много, стереть много(а стирать, по словам Скакова, медленнее, чем читать) и записать снова много. Вот тут и +1 к WA. Данное, реально записанное на флеш-память(Данные, отправленные на запись кистом. Все что дело связано долгие. Поэтому контроллер ищбы на заранее оговоренные иешей и ищет информацию о физическом нахождении части файла в таблицах трансляций и пометает старый блок как устаревший, но не пометает. По мере жизни SSD все больше таких простаивающих блоков получается. Что делать? Идея(иногда от трупа, пока ищба не выйдут(пока нет знака на запись и чтение), но тогда, может быть, делем лишнюю работу, так как при стирании ненужных блоков мы должны перемещать данные, которые потенциально(иногда) удаляются(если конечно что-то не перемещать ограниченно). Зато последующая запись на очищенный сектор будет быстрее. Это фоновая оиства. А можем делать непосредственно во время записи новые нужные данные(новая оиства). В итоге стараемся комбинировать два способа очистки, оиства ущерб от обком. По поводу стандартных значений - не понятно. Паша говорит что то про WA=1, 2... < 1, < 1, если контроллер умеет считать данные.	он чего говорил, что какие то бенчмарки тестировали на наборе 0 и было быстрое зашифрование. Разве Write Amplification это не про отношение количества информации, которое реально записывается на диск к количеству информации, которое должно быть записано? И Write Amplification > 1 когда, например, мы записываем что-то на область SSD, которая занята. Тогда мы стираем большой блок, а потом восстанавливаем его с измененным сектором и получаем, что надо было записать 1 сектор, а мы записали целый блок сектора, а < 1 когда у нас диск имеет в скане данных и вместо большого объема пишет меньший, скажем	
8	Почему размер сектора жесткого диска равно 4096?	переход на секторы размером 4 КБ позволил увеличивать плотность записи данных и емость жестких дисков, а также повышать надежность исправления ошибок.		
9	Почему размер сектора жесткого диска равно 4096?	Рассмотрим первый случай: 8 секторов по 512Б. У каждого из них существует M (Тбайт/Т) служебный иешей, с помощью которой можно в каждый блок восстанавливать до M / 2 ошибок. Во втором случае: 1 сектор 409Б и 8 М служебный иешей. Заметим, что суммарно по 8 блокам в первом случае можно исправить 4М ошибок, а также во втором случае суммарно 4М ошибок. Но так как сектор в первом случае по 512Б, то стиранию в одном секторе можно исправить не более M / 2 ошибок, а значит, что если в каком-нибудь секторе возникнет, например, 2М ошибок, то их исправить мы уже не сможем. Размер в 409Б связан с тем, что у файловой системы кластер в 4КБ.	<- в нескольких подробные ты конспект про ФС?	
10	Почему запись на диск происходит посекторно, а не побайтово?	Посекторно оборудование некачественно, то иногда некоторые биты могут записываться неправильно, поэтому существуют ECC и CRC, которые фиксируют ошибки. Сначала, что сделать: транза для сектора была бурной, а транза будет заново, имень нехот служебной информации, чем транза была инась еще дохера байткой для того, чтобы гарантировать корректность. Ну и у устройства DOR-памяти ишо, что перезапис сразу несколько байткой быстрее, чем один (эта одна из модификаций)		
11	Когда чтение диска не происходит?	Мы захотели прочитать, то есть прочитали иешбу то, что записали, потом сделали ECC, оно типично исправило ошибки, потом переписали CRC и сравнили с тем, которое есть, если оно совпало, то мы прочитали, если нет то ничего этого нам не отдает (слушайте диалог записку, транзе 2, 2 и немного 1,6 тайминг 34,29)	Помните за вопрос пайе у тебе есть диск(магнитный или оптический), ты хочешь его почитать, но он не дает тебе данные, такое бывает когда они там ошибочные на чтение был пример с телесом, что тот вначале (когда мало ошибок) фиксировал их, а потом резко выключалось, так это потому что он больше не может их фиксировать, а ошибочную информацию не выдает, такие и в дисках	

10	Главные отличия магнитных и оптических носителей	Способом хранения информации: магнитный диск хранит информацию с помощью магнитных полей, а оптический - с помощью ячеек на дорожках	
11	В чем схожи оптические накопители и flash-память	Нет так тасовой операции перезаписи: ты вначале все стираешь полностью (возможно, тебе придется записать куда-то данные, которые ты не хочешь терять), и потом на место записываешь то, что считываешь	кажется неполным ответ... (сделайте пометку правее если что-то изменилось позже)
12	Чему равен сектор при записи аудио?	2 352 байта, т.е. при записи аудио нам все-равно, что некоторые биты будут неправильными (человеческое ухо не поймет разницы), поэтому мы отдаем привычное место ECC и CRC под звук. Но такое есть специальные требования к самому аудио файлу: стерео, 16 битное кодирование, 44 100 Гц частота	
13	Есть два диска в черных коробках, один - магнитный, другой - оптический, как определить, где какой?	Есть умное решение, как например пометить их в микроволновку. Оптический диск будет жариться с потрескиванием, и слабым изворами (которые в коробке то и не увидишь). А вот магнитный выдлет неплохую порцию или вообще бегит от этой микроволновки хорошее решение. Я (оно одобрено?) Сергеем???? (с) Автор одобренное подраисовка бля в зеленых) Еще есть версия поговорить использовать магнит (но нуен стель мощный магнит), а также погнати (но я, что то, не верю, что это рабочий способ)	Стоит сравнить, как изменяется скорость считывания. Далее нужно понимать, что скорость считывания сектора линейно зависит от (длины сектора) / (длины окружности, на котором находится сектор) при фиксированной угловой скорости. Оптический диск: однослойный - скорость будет равномерно только увеличиваться двуслойный - скорость равномерно увеличивается, затем равномерно уменьшается. Сказано это тем, что у нас две дорожки (в однослойном диске), начинающиеся в центре диска, в которой длины секторов одинаковы. Магнитный диск: скорость считывания уменьшается "скачками", т.е. считываем, начиная с внешней дорожки, и дорожки разбиты на "группы", где каждая дорожка разбита на одно и тоже количество секторов.
14	RAID 0 1, 0+1, 4, 5, 6, 10 и +, отличия от остальных + чем лучше RAID-Солонки	RAID 0 - Redundant Array of Inexpensive Disks - Обычный массив недорогих дисков. Альтернатива SLED - Single Large Expensive Disk. RAID 0 - просто несколько дисков, данные распределяются по всем. Не использует ни бит четности, ни какой-либо ECC. Вообще-то так себе RAID, т.е. никакой избыточности в нем нет. Поддерживает диски разного объема, но объем памяти, который добавляет в массив каждый диск будет равен объему наименьшего из дисков. Теоретически, если у нас есть n дисков, то мы улучшим скорость передачи в n раз. RAID 1 - каждый диск - точная копия первого. Массив может быть настолько большим по памяти, насколько большой у нас самый маленький диск. Не использует бит четности. Работает Солонкина и всё такое. Повышает скорость чтения, но ухудшает скорость записи. RAID 10 - (n/2) дисков представляет собой RAID 0 массива, а n-ный диск хранит в себе бит четности всего массива. Неплохая скорость чтения, но с записью всё плохо. Кроме того, диск с битам четности используется сильнее чаще остальных и потенциально отвалится первым. RAID 5 - теперь блок с битам четности распределен по всем дискам. RAID 6 - теперь у каждого диска есть еще и блок с кодами Рада Солонкина, что повышает надежность. RAID 0+1 - RAID 0 массив из RAID 1 массивов RAID 1+0 - RAID 0 массив из RAID 1 массивов Коды Рада Солонкина позволяют при добавлении к сообщению длины N доли информации длины M исправить M/2 ошибок из всех (N+M) битов. Ну потому что менять состояние транзистора (особенно если это MLC или TLC) гораздо дороже, т.е. надо точнее падавать, например, чем поменять состояние с помощью магнитного поля. На SSD проблема решается тем образом, что используют докера плавок, и мы получаем какую-то пародию на параллельность	Таненбаум, страница 115 а на вычитайте еще и картиночки красивые есть
15	Почему флешка работает медленнее HDD?	Ну потому что менять состояние транзистора (особенно если это MLC или TLC) гораздо дороже, т.е. надо точнее падавать, например, чем поменять состояние с помощью магнитного поля. На SSD проблема решается тем образом, что используют докера плавок, и мы получаем какую-то пародию на параллельность	Возможно потому что во флешку не заложить умный контроллер
16	Как лучше подкармливать оптический диск: по дуге или вдоль радиуса	Если подкармливать диск по дуге, то часть данных пропадет и не восстановится. Если вдоль радиуса, то код коррекции должен справиться.	
17	Что покрывает код: ECC или CRC?	ECC покрывает CRC, поскольку ECC может поправлять ошибки, а CRC лишь чекат, что ошибок нет.	
18	Зачем сектор во флеш	Для коммек: сначала делают CRC, затем ECC, CRC - хэш-функция, ECC - код коррекции.	> Это унаследованный от жестких дисков интерфейс <- и ч-м не новый не сделали? > Хм, возможно чтобы для системы выглядело как HDD, если контроллер памяти на процессоре, ему проще пришлось бы объяснить, почему пропали сектора и как этим пользоваться, а так он ничего не записывает. Не забывайте, что переиспользование данных происходит тем же образом: "запомнить что было в секторе до этого, стереть всё, записать на место, не теряя при этом нужных старых данных и добавив новых", мб поэтому они и нужны во флеш. Н. Странно не поспешить, а побоясь, блок из памяти секторов, в террии ему сектора и не особо нужны для стирания и перезаписи. Еще вариант, сектора нужны, потому что ECC и CRC так не позволяют различать. Мб запишем несколько версий. И давай же каждую новую ревизию выделять отдельным сектором. И ты уверен(о), что запись не происходит не поспешно (см. 96 строку) > начинай @96b12 Н. Почему читать по одной ячейке развозит ресурс памяти? Мы же все равно и так и так читаем блок ячеек. А без сектора мы бы вообще читали не весь сектор (512 байт, например), а только то, что нам нужно. > Не побоясь же чтение делать... Понятно оно только для того, чтобы оптимизировать количество и эффективность операций чтения-изменения-стирания-записи. Работает оно примерно как с оперативкой (там же интерферсные семы, похиение) После команды чтения сектор копируется в внутренний буфер из SRAM. Если данные приходят команда на запись — операция изменения остается в этом буфере. Потом сбрасывается. Если бы чтение было побоясь, то при последовательном чтении (да и не только) интерес был бы забит командами. Т.е. там должны быть какие-то минимальные единицы информации для обмена с устройством, это как раз и есть сектор. Почему сеченя решил считать, 512 байт — вот ту может быть совместность (хотя не жейству) - @96b1wmt
<b>2.1 Архитектура фон Неймана и её альтернативы</b>			
1	Рассказать основные принципы.	1. Следует использовать децентрализованную систему считывания. 2. Программный принцип управления, т.е. обеспечивает универсальность разработки. 3. Адресность памяти: у каждой ячейки памяти есть свой номер и любой ячейке должен быть быстрый доступ. 4. Параллельность выполнения команд. 5. Однородность памяти, т.е. в одной памяти хранятся и данные и команды	
2	Рассказать про гарвардскую систему. Сказать что модифицированная есть.	Гарвардская архитектура - те же принципы, что и в фон Неймане, кроме последнего. В Гарвардской архитектуре разные каналы и физические устройства хранения для данных и инструкций. Про модифицированную версию: есть несколько вариантов, все они используют принцип "а давайте немного ослабим разделение между данными и инструкциями". Сказав рассказав про эй-сис-те, у нас разные каналы для данных и инструкций, но память общая. Собственно программист может иногда и не узнать, что у него модифицированная гарвардская, но выигрыш в скорости он получит	
3	Привести альтернативы каждому принципу фон Неймана	1. Дестинация(?) система считывания 2. Аппаратное управление 3. Машинный Тьюринг 4. Суперскаляр OOO 5. Неоднородная память (например Гарвардская архитектура)	Чем хуже? 1) 2) программное управление универсальнее, можно прогать разные алгоритмы, а аппаратное это делать сложнее под определенный апотс 3) у машины тьюринга относительная память, можете обматываться вами 4) суперскаляр может оптимизировать выполнение параллельных действий - по времени выигрыш, но железа сложнее и дороже 5) однородность проще, можно dimensionally строить код, но медленнее и есть баги в безопасности
4	Принципы фон Неймана. Почему десичная запись (первое машины были 10-ные)	Сначала пытались сделать 10-ные, но затем поняли что: десичную систему проще реализовать (достаточно записать/разложить ячейку памяти) довольно эффективная скорость вычисления (лучше только троичная система считывания, но она хуже исследована в сложение)	Если Паша строит, почему не извясился 3-чей, скажите, что использование 2-чей сложнее исторически, а на 3-нюю никто перекардять не будет. И 2-чей хорошо изучена и исследована, а так же не выгодно писать софт, придумывать стандарты и тд, под новые 3-чей компы
5	В чем разница между фон Неймановской архитектурой и Гарвардской с точки зрения программистов?	В Гарвардской архитектуре принципиально невозможно осуществить операцию записи в память программ, что исключает возможность случайного нарушения управляющей программы в случае ошибки программы при работе с данными или атаке третьих лиц (http://dmitry.korotkiy.ru/)	
6	Чем лучше: Аппаратное или программное управление? (любимый вопрос Саши)	Аппаратное невозможно исправить ошибки в уже выпущенной железе. Отказывай их и делать новое дороже. Можно в документацию конечно прописать что соре за ошибку, но это не очень хорошо. < в сложение хануть, потому что железа умеет делать, только то что вышло при производстве < жрет меньше энергии. Иногда вместе с программными вставками аппаратные блоки, реализующие какую-нибудь штуку (например видео декодер), а целью оптимизации энергопотребления и тепловыделения. < с другой стороны заточено под определенный стандарт, например, если видео в каком-то не поддерживаем формате, то декодер может не справиться	Программное - возможность исправления ошибок(выпустить апдейт и готово) < просто в разработке. Берем готовый универсальный вычислитель, который уже был сделан, и раздуем кучу. < с другой стороны железа будет дороже < для специализированного железа лучше сделать аппаратное, потому что стоит будет дешевле. Выбравший такие модели аппаратное, чем сложнее дороже железа, в котором будет использоваться лишь одна функция(все зависит от объема производства, в маркетинге упрощает кароч)
7	Есть два устройства: одно на фон Неймане, другое на Гарвардской. Как отличить? (включать нельзя)	На кристалле у Гарвардской будут две шины: для данных и для команд. А у фон Неймана только одна.	
8	Есть два устройства: одно на фон Неймане, другое на Гарвардской. Как отличить? (открывать коробку нельзя)	Программно: Однородная память: ячейка может быть или данными или кодом, в зависимости от того как мы с ней общаемся. Есть читать, записывать, то воспринимаем как данные. Если надо то как команду. Неоднородная память: у нас отдельные указатели для памяти и команд. Адресное пространство разделено на участки, каков-то область адресов отвечает за память, а другая за команды. К каждой области мы можем обращаться только по своему. Т.е. к памяти данные нельзя обращаться как к командам.	
9	Есть два устройства: одно на программном управлении, другое на аппаратном. Как отличить? (ниг, не посмотрев на кристалл - он может быть одинаковым, только у аппаратного отрублена программируемая часть)	Ну казалось бы они будут отличны внешне, когда разберем. Этот вариант, ему конечно тоже можно вынуть. Но. Вниманию. Они могут быть внутри абсолютно одинаково устройены, но у аппаратного может быть аппаратно запрограммирована реакция. Можно еще вынуть про то, что у программного будет версия прошивки, и соответственно его можно обновить. Можно вынуть, посмотреть в документацию(этоzero тупо конечно, но на уровне броска неплохо)	
10	Есть ли отличие между блоком данных для команд и данных в Гарвардской архитектуре?	В гарвардской архитектуре характеристика устройств памяти для инструкций и памяти для данных не обязательно должны быть одинаковыми. В частности, шина доступа, тактовая частота, технология реализации и структура адресов памяти могут различаться.	Длина слова, тайминги, технология реализации двух блоков памяти могут быть разными. В некоторых системах программа исполнения записана в ROM'e, т.е. нет необходимости ее менять, а в блок данных может быть перезаписываемым. Где-то, например, требуется больше памяти для команд => вместимость блока и размер адресов для команд и данных могут быть разными.
<b>2.2 Архитектура набора команд (ISA) и микроархитектура</b>			

		<p>ISA - Instruction Set Architecture (Архитектура Набора Команд).</p> <p>Без ISA и без наш код, нашу программу, мы можем точно определить, как она будет выглядеть в оптимизированном виде. ISA определяет:</p> <p>1) Архитектуру памяти. Стек, аккумулятор, Reg-Reg, Reg-Mem, Mem-Mem</p> <p>2) Набор команд (доступных программисту, пишущему на машинном языке)</p> <p>3) Количество регистров на кристалле. (CISC (complex instruction set computer) - обычно до 32, RISC (reduced instruction set computer) - обычно больше 32)</p> <p>4) Обработка ошибок. (Протокол обработки ошибок, типично не как делает наш код в случае если случилось говно. Конкретно см. внизу явочки)</p> <p>5) Рамочная адресация. (Видео, адрес памяти)</p> <p>6) Доступные типы данных. (целые, с плавающей точкой, знаковые, беззнаковые, векторные)</p> <p>7) Протокол взаимодействия с внешними устройствами ввода и вывода</p> <p>8) Режимы адресации:</p> <p>а) непосредственная адресация - хранение в адресной части самого операнда, а не его адреса или какой-либо другой инфы. таким образом можно работать только с константами.</p> <p>б) прямая адресация - хранение полного адреса операнда. команда всегда имеет доступ только к одному и тому же адресу памяти, значения в нем могут меняться, а вот сам адрес - нет. таким образом прямая адресация может использоваться только для доступа к глобальным переменным.</p> <p>в) регистровая адресация - запоминает прямую, только в данном случае вместо явочки памяти указывается регистр. благодаря быстрому доступу к регистрам и их коротким адресам это самый распространенный режим адресации.</p> <p>много компьютеров определяет также часто встречающиеся переменные (например, параметр цикла for) и помещают их в регистры (ее оптимиза).</p> <p>г) косвенная регистровая адресация - явочный операнд берётся из памяти или отправляется в память, но его адрес не фиксируется в команде, а содержится в регистре. если адрес используется таким образом, он называется указателем. преимущество - можно обращаться к памяти, не имея полного адреса. можно использовать разные слова памяти, меняя лишь значение в регистре. (посмотреть на стр.407 Андрюхи Таненбаума, там показано, как это может упростить жизнь).</p> <p>д) индексная адресация - обращение к памяти по регистру и константе смещения (от себя) - полагаю, что это напугивает обращение к явочке массива : ух азвель на начало + индекс смещения). Пример использования данного режима адресации см. стр.409 Таненбаума)</p> <p>е) относительная индексная адресация - адрес вычисляется путем суммирования значенй явочки двух регистров и смещения, один из регистров - баз, а другой - индекс. (MCU RA, PC + RS)</p> <p>ж) "начинанием упарываться"</p> <p>з) стековая адресация - максимально короткие команды, настолько сука короткие, что в них нет адреса. это возможно при наличии стека. при работе со стеком обычно используется постфиксная запись.</p> <p>ВАЖНО!</p> <p>Хорошая архитектура должна определять набор команд, которые можно эффективно реализовать не только в современной, но и в будущей технике.</p> <p>ОБРАБОТКА ОШИБОК</p> <p>Следующее предложение означает, что если программа выполняет код операции, который не определен, он должен вызывать системное прерывание, а не просто игнорироваться.</p> <p>Выполнение зарезервированного кода операции должно вызывать системное прерывание.</p> <p>Может быть и альтернативный подход.</p> <p>Результат выполнения зарезервированного кода операции определяется реализацией.</p> <p>Это значит, что разработчик компилятора не может рассчитывать на какое-то конкретное поведение, у разработчика возникает свобода выбора.</p>	
1	Что такое ISA? для чего она? как она определяется? Что в себе содержит? Описание происходящего при ошибках (?)	<p>После появления ISA программы, написанные для одного компьютера будут работать на любом другом, с такой же ISA. Вне зависимости от аппаратной реализации этой ISA. Т.е. возникает некоторая "гиростоянка" (абстракция) между аппаратными и программными обеспечениями</p> <p>Обратная совместимость - способность новой машины выполнять старые программы без изменений.</p> <p>Новые компьютеры дополняют старый набор команд &lt;==&gt; один и тот же код работает на новом, если на старом работ.</p> <p>Все компьютеры различного типа лишь внутренним устройством.</p> <p>Кросс-обратная совместимость: любые языки любого производителя имеющие одну ту же ISA могут выполнять программы, написанные для этой ISA</p>	про имеют одинаковый набор не очень грамотно. (Fixed)
2	Что стало лучше после появления ISA?	После появления ISA программы, написанные для одного компьютера будут работать на любом другом, с такой же ISA. Вне зависимости от аппаратной реализации этой ISA. Т.е. возникает некоторая "гиростоянка" (абстракция) между аппаратными и программными обеспечениями	
3	Что делает ISA для того, чтобы прог козались на всех железах?	Обратная совместимость - способность новой машины выполнять старые программы без изменений. <p>Новые компьютеры дополняют старый набор команд &lt;==&gt; один и тот же код работает на новом, если на старом работ.</p> <p>Все компьютеры различного типа лишь внутренним устройством.</p> <p>Кросс-обратная совместимость: любые языки любого производителя имеющие одну ту же ISA могут выполнять программы, написанные для этой ISA</p>	
4	Что такое Микрорхитектура и "с чем ее едят"?	В компьютерной инженерии микрорхитектура такое называемая организацией компьютера — это способ, которым данные архитектура набора команд (ISA) реализована в процессоре. Каждая ISA может быть реализована с помощью различных микрорхитектур. Например, Pentium и Intel Core - разные микрорхитектуры, которые поддерживают одну и ту же ISA (x86). Или System300. Глубоко, сложение 32-х битного числа может быть реализовано 32-х битным сумматором, а может — с помощью двух 16-битных. Haverel, Sandy Bridge, Nehalem, Skylake - x86 (еще Pentium и Intel Core)	
5	Приведите примеры Микрорхитектуры.	Со стороны ARM могу привести пример вых0 (vml4 - vml B)	
6	Примеры разных ISA (хватит двух) и разных микрорхитектур с одной ISA	ISA: MIPS, ARM, x86. Вытекли из Intel Core - разные микрорхитектуры, но одна ISA (x86) System300 - чудесная IBM/моякая линейка времен динозавров. Все модели поддерживали одну и ту же ISA, но реализована она была по-разному. (у дешевых моделей железо попроще, работали они, соответственно, дольше и в памяти имели меньше)	
7	Различия CISC и RISC	<p>CISC:</p> <ul style="list-style-type: none"> <li>- переменная длина команд (1-15 байт) (частные команды занимают меньше места)</li> <li>- сложные команды (переменной длины - сложно декодировать)</li> <li>- Reg-Mem или Mem-Mem</li> <li>- Арифметические команды выполняются целиком за один раз</li> <li>- Сложно параллелизуются</li> <li>- мало регистров. Регистры специализированные.</li> </ul> <p>RISC:</p> <ul style="list-style-type: none"> <li>- быстрее, чем CISC. Т.к. места под часто выполняемые команды на кристалле больше, длина инструкций фиксированная (проще декодировать)</li> <li>- больше регистров и их общего назначения</li> <li>- простая адресация и опко. Многие только смещение взять</li> </ul> <p>- сложная система обращения к памяти</p> <p>Про что это? О том, что можно писать в квадратных скобках в команде.</p> <p>Возможно сложное обращение к памяти, типа умножить на что-то регистр, прибавить еще что-то</p>	<p>в риске простая адресация к памяти?? у меня в конспекте ровно наоборот &lt;==&gt; соря, твой конспект говно</p> <p>А можно услышать аргумент, почему она простая и почему она сложнее?</p> <p>Чет я в этот пункт не вдуваю</p> <p>Если вопрос в целом, то ответ: не за декодера, ему надо декодировать не 1 команду, а 3-4 (у нас же параллельные конвейеры), а определить где начало 3 или 4 команды, когда у них не фиксированная длина это классическое разделение, при фиксированной можно просто оптимизировать 4 декодера и нормально дождаться &lt;==&gt; опко, в RISC, деление адресации к памяти, потому что там нет чудесного режима адресации Scale или у, например, x86 с адресами вида 100 + 100 + 4*for, в RISC, придется сначала вручную посчитать эту прелесть, а потом уже обращаться к памяти. в CISC так адресоваться к явочке можно одной командой</p>
8	Что быстрее: стек, аккумулятор, reg-reg, reg-mem, mem-mem?	reg-reg Ну, смотря, самой операции сильно больше, чем регистров, потому и адреса в ней длиннее. <p>Но за этого команда для "mem будет занимать больше места, даже если оба аргумента — регистры, потому сам код больше и занимает больше места, а тут в числе в явочке.</p> <p>А если код забит бесполезными данными, то приходится обращаться к операции, которая сильно медленнее.</p> <p>Чем компактнее команды — тем они лучше кэшируются, тем быстрее всё и работает.</p>	
<b>2.3 Конвейерная архитектура, конвейер MIPS</b>			
<p>На конвейере к примеру требуется 7 тактов, чтобы выполнить 1 команду. Но ведь можно заставить 3 такта (или вообще сделать супер железу и будет нужен 1 такт). Получается конвейер ухет?</p> <p>Почему бы не сделать процессор с лучшей явочкой на конвейере, как например Pentium 4, у которого 31.</p> <p>2. Считаем, что возможно так оптимизировать, что каждая малая деталь будет работать достаточно быстро, что это получится так же по скорости, как и нормальный конвейер.</p>			
3	Что выполняется на каждой стадии MIPS конвейера на примерах команд JMP, LD, SUB	<p>Конвейер лучше, потому что:</p> <ol style="list-style-type: none"> <li>1. длина такта у конвейера короче</li> <li>2. больше полезного деления за единицу времени (выше IPC Instruction Per Cycle) т.е. по максимуму используем железо</li> <li>3. цепочку конвейера выше latency, но сильно меньше throughput. Мы как бы жереем временем исполнения одной команды от начала и до конца, чтобы получить команду каждый такт.</li> </ol> <p>Вариант 1: Pentium (просто из лиза))0</p> <p>Ну, допустим у нас 5 конвейеров там был колоссальный такт с очень высокой частотой частоты (если 3 и 4 пентийму запускать на одинаковой частоте, то 4 проигрывает. Собственно поэтому после него Интел обратно откатился на архитектуру тригетто)</p> <p>Вариант 2: При любом JMP будет очень неэффективно использоваться конвейер, т.к. придется отплатиться на очень много стадий назад.</p> <p>Вариант 3 (допустимый SuperScalor): хазаров станет так много, что станет больно и дорого, упадет IPC. Будет жрать энергию как тварь.</p> <p>JMP:</p> <p>IF - прочитали команду из памяти</p> <p>ID - поняли, что это JMP, поменяли значение в Program Counter (Сюда называет эту штуку Instruction Pointer как в x86)</p> <p>EX - NOP</p> <p>MEM - NOP</p> <p>WB - NOP</p> <p>LD:</p> <p>IF - прочитали команду из памяти</p> <p>ID - декодировали</p> <p>EX - посчитали адрес в памяти (сложили регистр и offset)</p> <p>MEM - прочитали</p> <p>WB - положили в регистр</p> <p>SUB:</p> <p>IF - прочитали команду из памяти</p> <p>ID - декодировали</p> <p>EX - посчитали</p> <p>MEM - NOP</p> <p>WB - положили результат в регистр</p>	
4	Может конвейер. Все стадии + пример. Все инструкции (приводятся по всему absolutely) (х2 раза было)	<p>1. IF - Instruction Fetch - считываем команду из памяти, передаем дальше</p> <p>2. ID - Instruction Decode - декодируем инструкцию, т.е. определяем что делаем, откуда берем данные и куда пишем. Считываем данные из регистров</p> <p>3. EX - Execute - выполняем команду либо считаем адрес в памяти</p> <p>4. MEM - пишем в память/читаем с памяти</p> <p>5. WB - Write-back - пишем в регистр</p> <p>Примеры смотри выше</p>	
5	Что будет, если каждую стадию конвейера разбить еще на подстадию и до какого момента мы сможем так получать выигрыш?	Будет выигрыш по тактовой частоте и сильно проигрывать по энергопотреблению. В терси можно уменьшить время одного такта до тех пор, пока мы сможем писать в регистры (которые у нас между стадиями) без явочки, и пока у нас ток от тактового генератора успевают поступать везде). На практике получим много хазаров и большое энергопотребление (см. выполнение 2-ой команды)	
6	За сколько минимально можем выполнить операции без конвейера и на нем?	Ну, допустим у нас 5 конвейеров 5 стадий. Latency (время выполнения самой первой команды) у нас 5 тактов, зато следующую команду мы получим уже на следующей (шестой) такт.	
7	Что дает сложение на этапе MEM MIPS-конвейера?	без конвейера мы можем выполнить одну команду за условие 3 такта, но и следующую мы получим еще через три "_,_/_/_/_"	
8	Откуда IF берет команду?	Обращается к памяти по адресу лежащему в IP. IP инициализируется числом при старте системы, с этого числа начинается BIOS	
9	Что такое IP физически?	Специальный регистр	
<b>2.4 Проблемы конвейера (hazards) и пути их решения.</b>			
1	Привести пример RAW хазарда + решение, Control + решение	<p>RAW:</p> <p>ADD R1 R2 R3</p> <p>SUB R4 R1 R5</p> <p>Как фиксировать Forwarding, "это не bar, а fian", "а давайте конвейер постоит, пока дождется, а потом начнем следующую инструкцию"</p> <p>Control:</p> <p>ADD</p> <p>JMP A</p> <p>SUB</p> <p>ADD ...</p> <p>ADD ...</p> <p>... - вот это тоже начнет считаться, а зазз</p> <p>A XOR</p> <p>Как фиксировать: исполнять JMP еще на стадии ID (ADD отвалится) + "это fian")00"</p>	тут есть подобно со слов Сидорова
2	Что такое структурные хазарды (нет, не предзадание переходов)?	Это когда у нас железо не может выполнять некоторые команды одновременно. Например если у нас одна шина и на инструкции, и на данные, то когда у нас выполняются LDSD на стадии MEM, то мы не можем считать новую команду на стадии IF, потому что один канал до памяти и мы не можем одновременно обращаться к памяти по одному каналу. <p>Как фиксировать:</p> <ol style="list-style-type: none"> <li>1. Дублировать железо. А именно, если у нас два канала до памяти, то такой баг не произойдет (граварская архитектура поэтому и быстрее). Она, а кши первого уровня обычно как раз и делится на команды и данные.</li> <li>2. сказать в документации "ну, типа, когда читаем/пишем знаем много времени"</li> <li>3. Опять же, ID помнит, что на IF ничего разумного не пришло в этот момент и имеет NOP</li> </ol>	(под редакцией @wily_liger для обсуждения пишите тут в комментариях или в тг)
3	Могут ли возникать WAR и WAW конфликты на MIPS?	NET. Почему? Для них вроде нужно изменение порядка выполнения команд, т.е. когда команды выполняются не последовательно, поэтому они возможны на SuperScale с OoO, но невозможны на MIPS и на суперскалере с OoO	
5	Как именно устроен forwarding? (видимо на уровне реализации, нуки тут)	Рисую предположить, что есть регистры, чисто чистого равно мощности множества (EX, MEM, WB), и после вычисления результат кладется в регистр, причем эти регистры устроены по принципу очереди и, возможно, когда ID видит, что дождется результат и вернуть его в нормальные регистры не получится, меняет адрес аргумента на адрес в этой той очереди.	Предполагаю что ID видит что нам нужна информация, которая будет доступна через forwarding. Например есть у нас регистры промежуточный после EX. Соединяем их обратно с EX. Тогда, когда нам в EX нужны те значения, которые еще не легли в основных регистрах, мы идём в промежуточный регистр после EX и берем значения оттуда. За это все отвечает ID, он умный (@wily_liger)
6	Что хидает NOP?	это может делать стадия ID. Или, как в MIPS - не делать. И тогда либо программист сам руками добавляет NOP, либо сам дурак	
7	Код Wai и Waw	в явочке таких хазардов нет (см вопрос 3)	
<b>2.5 Суперскалярная и VLIW архитектуры.</b>			

	VLIW: + Программист должен сам заботиться о оптимальном заполнении конвейера + Поддержка старых архитектур становится невозможной (например, добавим еще один конвейер, и все, исполнители что куда раскинулись) + Хороший компилятор написать ну очень сложно + Относительно простую(ну в сравнении с) архитектуру + Меньше места занимает + Проще и дешевле + Так как работает с кардой компилятор он может оптимизировать лучше, чем Scheduler, поскольку у компилятора в распоряжении сильно больше времени и памяти	Суперскаларная: Если планировщик умный, то: - Runtime оптимизирует выполнение данных кода и чего нибудь еще + Оптимальное использование ресурсов(если как в шахматах считать на пару ходов вперед) + NOP-ы переигреть сам планировщик > больше команд влезают в кэш. + ХАХАААААААААА! Много Хахадов. (если ОО!) - Планировка дорожек по энергетике(на него полезного не считают) - Планировка ограничений по памяти, поэтому может работать хуже чем компилятор (которому доступна вся оператива)	- последний плос VLIW-a потенциально некорректен, он казал что у алива есть больше времени и памяти. - у суперскалера и VLIW-a одинаковый набор хахадов: Ray Waw War. Это их общий минус, и его надо убирать
1 Почему минусы VLIW и суперскаларной архитектуры			
2 Пример работы со VLIW	код для суперскалара: ADD R1 R1 R2 ADD R3 R3 R4 ADD R5 R5 R6 ADD R7 R7 R8 ADD R1 R1 R3 ADD R5 R5 R7 ADD R1 R1 R5		
3 Привести программу для VLIW и суперскалара, которая считает сумму чисел в восьми регистрах	ну и можно выгрузить результат из R1 куда-то в память командой LD, для алива тот же прикол с параллельностью, только NOP надо самому писать (именно потому и назвали суперскалар, кеЮД)		
4 Каким образом происходит планировка при компиляции, почему это плохо?	Это плохо, потому что у компилятора сильно больше времени и ресурсов на планирование.		
5 Планы динамического планирования? Или пример, когда динамическое число (он же может тоже спросить?)	Планировщик знает про железу сильно больше, чем компилятор. А еще он может подстраиваться под ситуацию (компилятор один раз скомпилит, и все). Например, при кэш-промахе планировщик знает, что ничего связанного с текущей командой запроса к памяти ему ближайшие 100 тактов не светит, и поэтому переставляется и считает что-то другое.	-Что значит "таким образом"? Просто компилятор берет и передает программу во VLIW-архитектуру? Или берет и просчитывает, как распределить по доступным конвейерам	
6 Может ли планировщик нарисовать обе архитектуры.	да		
7 Попросили написать код для VLIW и Superscalar, который выполнит по равнению, но выдаст один результат			
8 Почему код команды более компактный в суперскаларной архитектуре	Потому что не храним NOP в коде.		
9 Что быстрее?	Superscalar т.к. код получается компактнее и лучше эскируется		
10 Код War Waw	WAR 1) MUL R1, R2, R3 2) XOR R6, R1, R5 3) ADD R4, R2, R3 4) SUB R5, R4, R5 Поскольку операция MUL очень тяжелая, а XOR зависит от MUL, то вначале выполнятся 3 и 4 операции, и в XOR мы будем уже козять новое значение аргумента  WAW 1) MUL R1, R2, R3 2) ADD R4, R2, R3 3) SUB R1, R4, R5 Вначале выполняются операции ADD и SUB, а потом в R1 запишется результат умножения	А тут точно нужны 3 и 2 команды в waw и war соответственно? Без них разве не будет хазада?  Я просто скопировал код Зайвы  WAW точно не будет	
<b>2.6 Процессоры: общего назначения/поточные, ядра/многоядерные процессорные системы, одновременная многопоточность (SMT, HT)</b>			
1 Чем smp лучше суперскалара?	Быстрее [0.5, 2] раза в среднем на 20-30%. Позволяет эффективнее козять железо. Утверждается что железо в любом случае используется НЕ хуже, чем суперскалар		
2 Зачем вообще делать поточный процессор?	Потому что это выгодно, когда у нас много независимых задач, которые можно легко разбить на треды, и нам все равно, когда посчитается одна тред, важно только насколько быстро посчитаются все треды.		
3 Чем отличается многопоточность от многопроцессорности?	У каждого процессора свой контроллер памяти. У ядра общая память(ну, исключая кэши и регистры)		
5 Почему задача, решаемые на потоковом процессоре очень хорошо параллелизуются. И это их свойство потоковые процессоры используют лучше, чем CPU (т.е. у потоковых процессоров меньше простоя ядер, но кадры из которых можно кодить и чуку тредов закрутить)	Задачи, решаемые на потоковом процессоре очень хорошо параллелизуются. И это их свойство потоковые процессоры используют лучше, чем CPU (т.е. у потоковых процессоров меньше простоя ядер, но кадры из которых можно кодить и чуку тредов закрутить)	Везде бы все эти вопросы имели почти один и тот же ответ	
6 Когда SMT бывает вреден?	если у одного thread еле влезал в кэш, а мы решили запустить не 1, а 2, то они будут плохо сосуществовать и всё будет работать медленнее		
7 что такое Hyper-threading? Чем отличается кадр с HT от ядра без него?	См. SMT(punkt 13), HT это патент Intel на технологию SMT	в HT ядра в 2 раза больше регистров, умный планировщик, за ядро нет	
8 что такое потоковый процессор? Почему мы используем их в видеокарте, а в обычном процессоре почти нет?	В обычных процессоре это невыгодно. Т.к. большую часть программы нельзя разбить на очень-очень-очень много тредов. Получается что запускаешь что-то с другим тредом /тредками на чипы ядра потока процессора очень сильно проиграывают по скорости. (использим, что один тред нельзя выполнить на разных ядрах)	Везде бы все эти вопросы имели почти один и тот же ответ	
9 Когда мы получаем преимущество при использовании SMT?	Может 1) Число виртуальных ядер больше числа физических(значе зачем оно надо) 2) Новая задача либо хороша параллелизуема, либо ее надо отчитать мого ядра и это можно делить на разных ядрах 3) Каждый thread должен быть настолько "не очень хорошим", чтобы он оставался сборочные конвейеры, которые планировщик будет забивать задачками от других threads - треды должны не подтягивать за ресурсы (кош)		
10 Необходимые условия для работы HT. Какой планировщик больше всего подходит на HT?	Не получится без многопоточности, нужно писать правильные программы. Выбывает тулой, т.е. почти за бесплатно получат независимые инструкции	вопрос 15 связан	
11 Какие технологии GPGPU есть? Что было до CUDA/OpenCL	Есть CUDA - платформа и али от NVIDIA. Работает только с их картами. Код выглядит красиво, потому что все магия с видеокарточкой прячется в библиотеку. Есть OpenCL - открытая платформа. Поддерживается зелеными карточками, но плохо (потому что а кто тогда будет использовать CUDA). Код выглядит как не такое, потому что магию приходится писать руками. Из этого брали и ресурсами давали, создавали текстуры из массивов данных и шейдеры считали на этих текстур новую. Потом эти новые текстуры считали массово выходных данных.		
12 Почему нельзя написать ось под видео?	Потому что у GPU нет "облагов" из исполнительных ядер, которая есть в CPU, например системы разделения памяти, т.е. кадры з ядре не могут делиться памятью между Thread'ами и не защищает память одного thread'a от модификации другими.	а смейт??? Или тут спрашивают про то, почему нельзя вывести проут?	
13 Что такое SMT?	SMT это когда у chips есть реальные ядра, а виртуальные. Так как планировщик ядра не может вызвать 100% забивки конвейеров, то делают такой фиш: делятте удавом число регистров, и запускат дава потока(треда) вычислений. А что, они же по умолчанию независимы, так что планировщику проще, и железа хватает		

Раскраска по авторитетности источника ->>>>	смело можно использовать (официальная документация и/или советовал СККВ)	Литература (лежит в папке литература на диске)				
	почти можно доверять (не оф документация, содержащая догадки, но нормальная, а еще нормальные конспекты)	Цифровая схемотехника	<- простое			
	большой частью адекватно (более-менее адекватные таблички прошлых лет)	Computer Architecture: a Quantitative Approach	<- крайне рекомендую приложения А, В и С			
	полуправда	Таненбаум	<- можно найти процентов 40 от билетов			
	не проверялось					
	другое					

Доки	Другое:	Статьи (английский - это судьба (ц))	Презентации	Технические документации	Видосики	<b>СТАТИСТИКА</b>
<u>Самая важная дока</u>	Симулятор схем	<b>ОЧЕНЬ ИНТЕРЕСНАЯ СТАТЬЯ</b>	Про MIPS конвейер (на английском, но с картиночками)	Описание ISA ARMv8-A	HDD	<b>y2015</b>
безбашные конспектики	Рисовалка схем	<b>БЕСКОНЕЧНАЯ МЫШКА</b>		Описание процессора ARM Cortex-A76 (например, на <a href="#">3.4 ядрах с 14CPI</a> )		<b>y2016</b>
дока препода 34-35 (есть картинки, но много багов)	DJ Zayka	про тайминги оперативки		DDR		<b>y2017</b>
конспекты Жени	представление целых чисел в компдактере	оч интересные 238 страниц про ассемблер и конвейеры различных цпу		Память		<b>y2018</b>
билеты by y2017	представление вещественных чисел в компдактере	Референс по AMD64 (на страницах 221-223 в pdf'ке есть подробная инфа о MOESI)				
еще билеты by y2017	курс на coursera, если проебали. пару про стек, reg-reg и kek-kek, тут почти 1 в 1	глава из книжки про триггеры и все такое				
		<a href="#">про счетчики</a>				





