

# Questions on Computer Architecture for students

## Caches/Memory - I

### Prerequisite questions

1. Почему в современных процессорах используется кэш для основной памяти?
2. Скажите, каким свойствам желательно, чтобы удовлетворяли запросы (access patterns) в память, чтобы кэш использовался наиболее эффективно?
3. Скажите, Вы знакомы с тем как устроены и как организуются кэши (cache organizations) в современных процессорах? К примеру, можете сказать чем отличаются direct-mapped cache (кэш прямого отображения) и set-associative кэши?
4. (Если проблемы с предыдущим вопросом) Можете рассказать в принципе, что Вы знаете об устройстве и работе кэшей в микропроцессорах?

### Main questions

1. Допустим, мы работаем с set-associative кэшем. (На какие группы бит разбивается адрес обращения в памяти?) Как адрес обращения в память обычно разбивается на группы бит, соответствующие полям set, tag и offset? (обязательно уточнить про то, где старшие биты, а где младшие биты в адресе)
2. Follow-up question: А почему именно так располагают поля в адресе? Что будет происходить, если к примеру поменять set и tag местами?

Answer: Адрес делится на битовые группы, соответствующие tag, set и offset. |tag|set|offset| от старшего бита адреса к младшему. Если В - размер кэш-блока в байтах, то размер поля offset будет  $\lceil \log_2(B) \rceil$ . Разделение происходит именно таким образом, чтобы кэш-линии с адресами А, А + В, А + 2 \* В и т.д. попадали в разные set-ы. В случае, если поменять set и tag местами на один set будут отображаться линейные блоки размером  $B * 2^{(tag\_size)}$ , что будет приводить к потере чувствительности к spatial locality, conflict misses, потере эффективности.

3. (Hard, need to show params on the screen) Предположим, что мы рассматриваем direct-mapped cache с политикой замещения LRU в микропроцессоре с 32-битной адресацией. Будем считать, что это L1 кэш данных. Размер кэш-линии в данном кэше составляет 64 байта. При этом, при разработке процессора что-то пошло не так, и поля set и tag поменялись местами в адресе обращения (т.е. |set|tag|offset| от старшего бита к младшего). Программист Вася написал следующий отрывок программы

```
(a) for (int i = 0; i < 100000; ++i)
    for (int j = 0; j < 1000; j += STEP)
        std::cout << a[j];

(b) auto a = new Bstruct[1000];
    auto end = a + 1000;
    for (int i = 0; i < 1000000; ++i)
        for (auto ptr = a; ptr != end; ptr++)
            std::cout << *(std::uint64_t*)ptr;
```

- (a) Можете привести пример какого-либо значения STEP,  
(b) Можете привести пример размера структуры Bstruct,  
при котором производительность процессора с данным кэшем с измененным порядком set и tag будет для этой программы лучше, чем со стандартным (неизменным)?

#### Help-1:

Давайте начнем с того, что определим принципиальный момент, почему производительность программы с измененным порядком может быть лучше

#### Help-2:

Давайте попробуем написать общую формулу сначала для данной ситуации, не привязываясь к конкретным числам и посмотрим, каких параметров нам не хватает

#### Help-3:

Сколько бит в адресе будет приходиться на offset, а сколько на set + tag?

Answer:

- В данном случае главный access pattern LRU-friendly [А, А + К, А + 2 \* К, ..., А + 1000 \* К].
- Кэш - direct-mapped, это значит, что кэш-линии, отображенные на один и тот же set будут гарантированно перезаписывать друг друга. Будем далее считать, что В - размер кэш-линии

- в байтах, #sets - число set-ов.
- В случае нормального расположения битовых групп в один и тот же set будут попадать кэш линии с адресами  $A$ ,  $A + \text{\#sets} * B$ ,  $A + 2 * \text{\#sets} * B$
- В случае измененного расположения битовых групп в разные set-ы будут попадать кэш-линии с адресами  $A$ ,  $A + 2^{\text{tag\_size}} * B$ ,  $A + 2 * 2^{\text{tag\_size}} * B$ .
- В таком случае должно быть  $2^{\text{tag\_size}} = \text{\#sets}$ . Так как адрес составляет 32 бита и 6 битов отводится на offset, то 26 битов приходится на tag + set. А значит  $\text{\#sets} = 2^{\text{tag\_size}} = 2^{13}$ .
- Следовательно,  $\text{STEP} = \text{sizeof}(\text{Bstruct}) = 2^{13} * 2^6 = 2^{19} = 524,288$  (байт) = (512 КБ)

4. (Вопрос, если на 1 вопрос был дан неправильный (и быстрый) ответ, но при этом OFFSET на правильном месте; может быть вместо 3)

Дан кэш direct-mapped с размером кэш линии в 32 байта. Что будет выведено на экран? Что будет происходить с кэшем для данной программы (предполагается, что в ответе на 1 вопрос set и tag были поставлены наоборот)

```
struct A {
    std::uint32_t my_array[3];
    std::uint8_t value;
};

int j = 0;
for (int i = 0; i < 1000000; ++i)
    for (j = 0; j < 30; j += sizeof(A))
        std::cout << a[j];

std::printf("%d\n", j);
```

Answer:  $30 * \text{sizeof}(A)$  + уточнение про padding. Если ответ 600, то спросить всегда ли так. Будет перетирание внутри одного set, в первом вопросе должно быть |tag|set|offset|

## Branch prediction - I

### Prerequisite questions

1. Можете сказать, что такое branch prediction (предсказание переходов) и почему его используют в современных процессорах?
2. Скажите, а если у нас в системе команд процессора нет инструкций условного перехода (conditional branches), нужен ли нам branch predictor?
3. Можете рассказать принцип работы bimodal branch predictor (2-bit saturating counter)?
4. Вы знаете, что такое local branch predictor и global branch predictor, и в чем между ними отличие?

### Main questions

1. Пусть рассматривается локальный (two-level adaptive predictor) branch predictor. Про него известно, что для каждого РС бранча мы храним полную локальную историю, по которой с помощью бимодального предиктора мы делаем предсказание. Считаем, что в начальном состоянии бимодальный предиктор находится в состоянии Strongly Not-taken. Размера бранч предиктора хватает, чтобы хранить информацию о всех бранчах, присутствующих в рассматриваемых программах.
- ```
for (int i = 0; i < 100; ++i)
    if (i % 2 == 0 || i % 7 == 0)
        std::cout << "Even\n";
```