

Лабораторная работа №2	Б10	2022
Моделирование схем в Verilog	Хорохорин Андрей Сергеевич	

Цель работы

Построение кэша и моделирование системы «процессов-кэш память» на языке описания Verilog.

Инструментарий

- Компилятор Icarus Verilog version 11.0 (stable)
- SystemVerilog стандарта IEEE1800-2012
- Лучший текстовый редактор Vim с использованием [языкового сервера](#) для SystemVerilog
- [gtkwave](#) как средство для просмотра .vcd файлов
- L^AT_EX с расширением XeLaTeX для написания отчёта

Задача

Имеется следующее определение глобальных переменных и функций:

```

1  #define M 64
2  #define N 60
3  #define K 32
4  int8 a[M][K];
5  int16 b[K][N];
6  int32 c[M][N];
7
8  void mmul() {
9      int8 *pa = a;
10     int32 *pc = c;
11     for (int y = 0; y < M; y++) {
12         for (int x = 0; x < N; x++) {
13             int16 *pb = b;
14             int32 s = 0;
15             for (int k = 0; k < K; k++) {
16                 s += pa[k] * pb[x];
17                 pb += N;
18             }
19             pc[x] = s;
20         }
21         pa += K;
22         pc += N;
23     }
24 }
```

- Сложение, инициализация переменных и переход на новую итерацию цикла, выход из функции занимают 1 такт. Умножение – 5 тактов. Обращение к памяти вида `pc[x]` считается за одну команду.

- Массивы последовательно хранятся в памяти, и первый из них начинается с 0.
- Все локальные переменные лежат в регистрах процессора.
- По моделируемой шине происходит только обмен данными (не командами).

Необходимо определить процент попаданий (число попаданий к общему числу обращений) для кэша и общее время (в тактах), затраченное на выполнение этой функции.

Вычисление недостающих параметров системы

Вычисление размера областей памяти

Везде далее будет считаться, что байт = 8 бит. Для начала выпишем все параметры кэша и памяти, которые были даны в условии:

- **CACHE_WAY** — 2
- **CACHE_TAG_SIZE** — 10 бит
- **CACHE_LINE_SIZE** — 16 байт
- **CACHE_LINE_COUNT** — 64 линии
- **MEM_SIZE** — 512 Кбайт

Так как всего в нашем распоряжении **CACHE_LINE_COUNT**(64) линии, каждая по **CACHE_LINE_SIZE**(16 байт), то из них легко вычисляется общий объём полезных данных, хранимых нашим кэшем.

$$\text{CACHE_SIZE} = 64 \cdot 16 \text{ байт} = 1024 \text{ байт} = 1 \text{ Кбайт}$$

Исходя из того, что всего необходимо проиндексировать 512 Кбайт памяти:

$$\text{CACHE_ADDR_SIZE} = \lfloor \log_2(512 \cdot 2^{10}) \rfloor = 19 \text{ бит}$$

Исходя из того, что всего у нас **CACHE_LINE_COUNT**(64) линий и все их нам необходимо разбить в блоки размера **CACHE_WAY**(2), то количество блоков

$$\text{CACHE_SETS_COUNT} = 64/2 = 32$$

Для индексации внутри каждого блока будет использоваться **CACHE_OFFSET_SIZE**, поэтому его размер зависит непосредственно от размера блока **CACHE_LINE_SIZE**(16 байт)

$$\text{CACHE_OFFSET_SIZE} = \log_2 16 = 4 \text{ бит}$$

Осталось посчитать размер индекса в наборе кэш-линий **CACHE_SET_SIZE**, он должен индексировать все блоки, которые мы имеем, а их количество нам известно **CACHE_SETS_COUNT**(32).

$$\text{CACHE_SET_SIZE} = \log_2(32) = 5 \text{ бит}$$

Заметим, что тег адреса, номер блока и сдвиг внутри блока должен однозначно задавать адрес ячейки в оперативной памяти, поэтому можно ещё раз проверить верность расчётов, проверив следующее равенство:

$$\text{CACHE_TAG_SZ} + \text{CACHE_SET_SZ} + \text{CACHE_OFFSET_SZ} = \text{CACHE_ADDR_SZ}$$
$$10 + 5 + 4 = 19$$

Видим, что всё в порядке — равенство верно.

Подводя итог, ниже перечислим все необходимые параметры системы:

- **MEM_SIZE** — 512 Кбайт
- **CACHE_SIZE** — 1 Кбайт
- **CACHE_LINE_SIZE** — 16 байт
- **CACHE_LINE_COUNT** — 64 линии
- **CACHE_WAY** — 2
- **CACHE_SETS_COUNT** — 32 блока
- **CACHE_TAG_SIZE** — 10 бит
- **CACHE_SET_SIZE** — 5 бит
- **CACHE_OFFSET_SIZE** — 4 бит
- **CACHE_ADDR_SIZE** — 19 бит

Вычисление размерности шин

Начнём с вычисления размера шины адреса между CPU и Cache. В нашей системе обмен данными по этой шине происходит за два такта. За первый необходимо успеть передать данные о теге и о номере блока, а за второй — offset внутри блока. Больше ни для чего эта шина не используется, а значит её размер можно сделать следующим:

$$\text{ADDR1_BUS_SZ} = \max(\text{CACHE_TAG_SZ} + \text{CACHE_SET_SZ}, \text{CACHE_OFFSET_SZ})$$
$$\text{ADDR1_BUS_SZ} = \max(10 + 5, 4) = 15 \text{ бит}$$

По шине A2 кэш передаёт интересующий его адрес памяти в контроллер памяти, причём часть адреса offset не передаётся, так как в нашем кэше мы делаем адреса начала кэш линий кратными их размеру. Так как данная шина ни для чего больше не используется, то её размер вычисляется следующим образом:

$$\text{ADDR2_BUS_SZ} = \text{CACHE_TAG_SZ} + \text{CACHE_SET_SZ} = 10 + 5 = 15 \text{ бит}$$

Вычисление размерностей шин команд это тривиальная задача, но для неё нам потребуется вспомнить, какие вообще команды(сигналы) могут передаваться по этим шинам.

Cpu → Cache 0 — C1_NOP 1 — C1_READ8 2 — C1_READ16 3 — C1_READ32 4 — C1_INVALIDATE_LINE 5 — C1_WRITE8 6 — C1_WRITE16 7 — C1_WRITE32	Cpu ← Cache 0 — C1_NOP 7 — C1_RESPONSE
Cache → Mem 0 — C2_NOP 2 — C2_READ_LINE 3 — C2_WRITE_LINE	Cache ← Mem 0 — C2_NOP 1 — C2_RESPONSE

Из таблицы хорошо видно, что для передачи необходимых сигналов между Сру и Cache достаточно 3 бит за такт, а для Cache и Mem — 2 бита за такт. Таким образом, имеем:

CTR1_BUS_SZ = 3 бит

CTR2_BUS_SZ = 2 бит

Аналитическое решение задачи

Попытаемся решить задачу аналитически. Для начала немного более строго сформулировать приведённую в задании систему. Далее будет считаться, что ни процессор ни кэш не обладают ни асинхронностью ни параллельностью. То есть когда процессор отправляет запрос кэшу, он только ждёт ответа и ничего не вычисляет. То же самое относится и к записи. Помимо этого, в ситуации, когда кэш ждёт ответа от памяти, он сначала примет все данные от памяти и только потом начнёт взаимодействие с процессором. То есть никогда не возникает ситуации, что кэш одновременно общается по шине и с процессором и с памятью. В ситуации, когда кэш сообщает памяти новые кэш линии он не ждёт никакого ответа от памяти. Кроме того, в конце выполнения не обязательно все кэш линии будут совпадать с памятью, в подсчёте времени работы не будет закладываться время на то, чтобы полностью синхронизировать кэш и память с точки зрения того, чтобы ни оказалось ни одной кэш линии с параметром `dirty=1`.

Для того, чтобы упростить себе задачу давайте разделим все такты на два вида: когда процессор просто ждёт ответа кэша и когда процессор что-то вычисляет. Подсчитывать количество попаданий/промахов вручную мне показалось сложным, поэтому был написан следующий код, который симулирует работу кэша по модулю того, что явно не хранит полезные данные и, соответственно, не возвращает их на запрос.

Тип `Word` используется для указания количества байт, используемых при `Read/Write` запросах. Класс `CacheLine` тоже ничего умного не делает, это просто удобный контейнер, представляющий собой кэш-линию, без хранящихся в ней данных.

```
typedef unsigned long long Addr;

enum Word { WORD=8, DWORD=16, QWORD=32 };

struct CacheLine {
    const int size;
    bool valid = 0;
    bool dirty = 0;
    Addr tag = -1;
    int last_call = 0;
    CacheLine& operator = (CacheLine oth) {
        assert(size == oth.size);
        tie(valid, dirty, tag, last_call) = tie(oth.valid, oth.dirty, oth.tag, oth.last_call);
        return *this;
    }
};
```

Ниже описан класс, представляющий собой один блок(set) из кэш-линий. Функции `read` и `write` условно эмулируют операцию нашей модели кэша, при этом `read` ещё возвращает есть ли запрашиваемый адрес в блоке или нет. Наверное, самой интересной функцией является `find_LRO`, которая ищет кандидата на вытеснение из кэша по следующей логике: если есть свободная линия(`valid=0`), то вытесняем её, иначе вытесняется линия, к которой условное время последнего запроса было минимальным. За условное время запроса берётся просто количество запросов совершенных до этого. Помимо этого, в классе вычисляется количество вытесняемый кэш линий, которые необходимо записать в память. Эта информация пригодится для работы других классов.

```
struct CacheSet {
private:
```

```

vector<CacheLine> lines;
const int line_size;
int total_mem_pushes = 0;
public:
CacheSet(int way, int line_size) : line_size(line_size) {
    lines = vector<CacheLine>(way, CacheLine{line_size});
}

bool exists(Addr tag) {
    return find_by_tag(tag) != lines.end();
}

void read(Addr tag, int call) {
    auto it = find_by_tag(tag);
    if (it != lines.end()) {
        it->last_call = call;
    } else {
        *(find_LRO()) = CacheLine {line_size, 1, 0, tag, call};
    }
}

void write(Addr tag, int call) {
    auto it = find_by_tag(tag);
    if (it != lines.end()) {
        it->last_call = call;
        it->dirty = 1;
    } else {
        total_mem_pushes += find_LRO()->dirty;
        *(find_LRO()) = CacheLine {line_size, 1, 1, tag, call};
    }
}

int getMemPushes() const {
    return total_mem_pushes;
}

private:
vector<CacheLine>::iterator find_by_tag(Addr tag) {
    return find_if(lines.begin(), lines.end(), [&] (const CacheLine& line) {
        return line.valid && line.tag == tag;
    });
}

vector<CacheLine>::iterator find_LRO() {
    return min_element(lines.begin(), lines.end(), [] (const CacheLine& a, const CacheLine& b) {
        return tie(a.valid, a.last_call) < tie(b.valid, b.last_call);
    });
}
};

```

Теперь можно перейти к самому сложному классу — Cache. Он занимается тем, что принимает запрос к памяти, определяет в каком он блоке выполняет его уже на уровне блока. Но помимо этого он считает интересную нам статистику — количество попаданий в кэш, количество промахов и прогнозируемое общее время работы данных операций на реальном кэше. Ещё он считает количество совершённых обращений к памяти в переменной calls, чтобы обеспечивать каждый запрос корректным условным временем.

Стоит отдельно объяснить почему время считается именно таким образом. В условии задания написано, что, в случае попадания, первый такт ответа кэша происходит через 6 тактов после первого такта запроса. Это значит, что если данные были переданы за один такт, то на седьмой такт процессор продолжит исполнять команды, получив ответ на запрос, который он послал ещё на первом такте. Итого на одном запросе процессор простаивает 7 тактов. Но данные от кэша могут передаваться не один такт, поэтому была сделана внутренняя функция `transfer_lag`, которая по порции данных определяет сколько она будет передаваться. Так как обе шины данных в нашей

системе имеют одинаковую пропускную способность, эта функция используется для аналогичных расчётов при передаче данных от памяти к кэшу.

В случае, когда на кэш поступает запрос записи всё происходит аналогично чтению за исключением одного момента: передачи данных от кэша к процессору не происходит, передаётся только команда-флаг, что всё отработало, но это всё равно занимает один дополнительный такт.

Помимо всего вышеперечисленного, необходимо помнить, что когда мы вытесняем из кэша линию, которая помечена `dirty=1`, то необходимо послать дополнительный запрос к памяти, эта поправка вычисляется непосредственно в методе `get_time`, используя насчитанные значения в каждой из `CacheLine`.

Также есть «волшебная» функция `split_addr`, в которой происходит много битовых операций, чтобы просто разбить адрес части на `tag`, `set`, `offset`.

```
struct Cache {
private:
    const int sets_cnt;
    const int way;
    const int line_size;           // byte
    const int data_bus_size;       // bits
    const int mem_size;           // byte

    int total_hits = 0;
    int total_misses = 0;
    int total_time = 0;
    int calls = 0;

    vector<CacheSet> sets;

    struct InnerAddr {
        Addr tag, set, offset;
    };
public:
    Cache(int sets_cnt, int way, int line_size, int data_bus_size, int mem_size) :
        sets_cnt(sets_cnt), way(way), line_size(line_size), data_bus_size(data_bus_size),
        mem_size(mem_size) {
        sets = vector<CacheSet>(sets_cnt, CacheSet(way, line_size));
    }

    void read(Word word, Addr addr) {
        calls++;
        auto [tag, set, offset] = split_addr(addr);
        if (sets[set].exists(tag)) {
            total_hits++;
            total_time += 6; // Cache lag
            total_time += transfer_lag(word) + 1; // Cache -> Cpu
                                   // 8-9
        } else {
            total_misses++;
            total_time += 4 + 100; // Cache + Mem lag
            total_time += transfer_lag(WORD * line_size) + 1; // Mem -> Cache
            total_time += transfer_lag(word); // Cache -> Cpu
                                   // 114-115
        }
        sets[set].read(tag, calls);
    };

    void write(Word word, Addr addr) {
        calls++;
        auto [tag, set, offset] = split_addr(addr);
        if (sets[set].exists(tag)) {
            total_hits++;
            total_time += 6; // Cache lag
            total_time += 1; // Cache->Cpu response
        } else {
            total_misses++;
            total_time += 4 + 100; // Cache + Mem lag
            total_time += transfer_lag(WORD * line_size); // Mem -> Cache
            total_time += 1; // Cache->Cpu response
        }
        sets[set].write(tag, calls);
    };
};
```

```

int get_hits() const {
    return total_hits;
}

int get_misses() const {
    return total_misses;
}

int get_time() const {
    int total_pushes = accumulate(sets.begin(), sets.end(), 0,
        [](int acc, const CacheSet& el) {
            return acc + el.getMemPushes();
        });
    return total_time + total_pushes * (100 + 1);
}

private:
InnerAddr split_addr(Addr addr) const {
    return {
        (addr & ((mem_size - 1) & ~(sets_cnt * line_size - 1))) / (sets_cnt * line_size),
        (addr & (sets_cnt - 1) * line_size) / line_size,
        addr & (line_size - 1)
    };
}

int transfer_lag(int data) const {
    return (data + data_bus_size - 1) / data_bus_size;
}
};

```

Теперь давайте направим всё реализованное нами «чудо» на решение задачи. Код ниже не симулирует пример из задачи полностью, а лишь обращается с запросами чтения/записи к тем же адресам, но для анализа времени работы кэша это и требуется. Класс PseudoAllocator просто помогает симулировать последовательное выделение памяти для массивов, возвращая адреса начала каждого из них.

```

struct PseudoAllocator {
    Addr last_allocated_addr = 0;

    Addr allocate(int n, int m, int elem_sz) {
        int res = last_allocated_addr;
        last_allocated_addr += n * m * elem_sz;
        return res;
    }
};

int main() {
    // 512Kb == [0x00000...0x7ffff]
    Cache cache(32, 2, 16, 16, 512 * 1024);
    PseudoAllocator alloc;

    const int M = 64;
    const int N = 60;
    const int K = 32;

    int8_t* a = (int8_t*) alloc.allocate(M, K, sizeof(int8_t)); // a[M][K];
    int16_t* b = (int16_t*) alloc.allocate(K, N, sizeof(int16_t)); // b[K][N];
    int32_t* c = (int32_t*) alloc.allocate(M, N, sizeof(int32_t)); // c[M][N];

    for (int y = 0; y < M; y++) {
        for (int x = 0; x < N; x++) {
            int16_t* pb = b;
            for (int k = 0; k < K; k++) {
                cache.read(WORD, (Addr) (a + k));
                cache.read(DWORD, (Addr) (pb + x));
                pb += N;
            }
            cache.write(QWORD, (Addr) (c + x));
        }
        a += K;
        c += N;
    }
}

```



```

cout << "HITS: " << cache.get_hits() << "\nMISSES: " << cache.get_misses() << "\nTOTAL TIME: " <<
cache.get_time() << endl;
cout << "RATE: " << 1.0 * (cache.get_hits()) / (cache.get_hits() + cache.get_misses()) << endl;

return 0;
}

```

Ниже приведён вывод данного кода:

HITS: 228080

MISSES: 21520

TOTAL TIME: 4274080

RATE: 0.913782

Теперь, когда мы знаем за сколько работает кэш, вычислим время работы всего остального.

- Один такт на выход из функции.
- Инициализация переменных, каждая из которых выполняется за один такт. Она явно выполняется в строчках 9, 10, 13, 14. И неявно в каждом из циклов. Будем считать, что переменные локальные для цикла выделяются заново на каждой итерации внешнего цикла. Тогда суммарное количество инициализаций вычисляется как: $3 + M + M \cdot N \cdot 3 = 11587$
- Подсчитаем количество тактов на сложение, включив в них такты на смену итерации цикла. $M \cdot 3 + M \cdot N + M \cdot N \cdot K \cdot 3 = 372672$
- Количество тактов на умножение: $M \cdot N \cdot K \cdot 5 = 614400$

Итого дополнительно имеем: $1 + 11587 + 372672 + 614400 = 998660$. Тогда итоговое прогнозируемое время работы равно: $998660 + 4274080 = 5272740$

Моделирование на языке Verilog

Структура проекта

Стоит начать с общей структуры модели и с её организации на файлы. Логически есть 3 модуля: процессор, кэш и память, причём процессор в нашем случае не был реализован явно, а лишь в тестового окружения на Verilog, который лишь симулирует работу настоящего процессора для кэша.

Пройдёмся по файлам:

- `parameters.sv` для описания глобальных констант и структур.
- `clock.sv` в котором задаётся модуль `Clock`, который используется для синхронизации схем и оценки времени их работы, так как внутри себя он считает количество тактов.

- `cache.sv` и `mem.sv` задают основные модули данного проекта — `Cache` и `Memory`. Отдельно стоит отметить, что в этих же файлах располагаются удобные интерфейсы доступа к ним: `CacheDriver`, `MemoryDriver`. Они необходимы для избежания дублирования кода, так как к кэшу и к памяти есть нужда обращаться из разных модулей, а писать больше одного раза код, который описывает протокол общения с кэшем очень непродуктивно.
- `mem_testbench.sv` и `cache_testbench.sv` хранят реализацию модулей Unit тестов для памяти и кэша соответственно. Помимо корректности, они проверяют ещё и время работы модулей в различных ситуациях, что оказалось чрезвычайно полезным.

Для подключения модулей используется директива **include**, работающая аналогично одноимённой директиве в C и имеющая те же проблемы: а именно множественный `include` одного и того-же модуля и, как следствие, ошибка компиляции. В Verilog эта проблема решается точно так как, как и в C, при помощи так называемых `include guards`.

```
`ifndef CPU_GUARD
`define CPU_GUARD
...
`endif
```

Inout шины

Поясним, как реализованы inout шины на Verilog. Для примера возьмём часть кода из файла `mem.sv`, в которой реализована шина адреса `cmd_w`.

```
module Memory
(
    inout wire[1:0] cmd_w
);

    bit owner = 0;
    logic[1:0] cmd;

    assign cmd_w = owner ? cmd : 2'bzz;
```

Для каждой шины в модуле заводится `bit` переменная обозначающая, владеет ли данный модуль шиной в текущий момент времени. В данном примере она названа `owner`. Помимо этого, каждый **inout** выход модуля имеет `logic` переменную внутри модуля, данном случае это `cmd`. Она необходима для того, чтобы взаимодействовать с шиной `cmd_w` как `output`, так как напрямую записать значение в неё не получится. Поэтому делается интересный трюк при помощи конструкции

assign `cmd_w = owner ? cmd : 2'bzz`, которая в конце каждого такта присваивает значение `cmd_w` либо `cmd`, если наш модуль должен выставлять значение на

шине, либо высокоимпедансное состояние для того, чтобы на проводках `cmd_w` принимались значения, которые выставляет другая сторона и не появлялось конфликтов сигналов, которые могут реальной схеме даже привести к короткому замыканию и натурально поджарить нашу схему.

Так как описанный выше трюк приходится осуществлять с любой шиной, работающей на ввод-вывод, то эти подробности реализации помещены в модули-драйверы, для удобства моделирования вышестоящих уровней.

Почему размер шин адреса не совпадает с минимальным достаточным?

Отдельно хочу отметить, что размер шин адреса `ADDR1_BUS_SIZE` и `ADDR2_BUS_SIZE` в моей реализации округлены до 16, для того, чтобы все шины адреса и данных единообразно задавались целым количеством байт, мне показалось это удобным программировать.

В действительности, если бы я писал код заново, лучшим решением было бы хранить все размеры в битах.

Функция skip

Почти в каждом модуле можно видеть следующую функцию, которая необходима для удобного программирования схемы.

```
task automatic skip(input longint ticks = 1);
    logic enter_clk = clk;
    while (ticks > 0) begin
        wait(clk != enter_clk);
        wait(clk == enter_clk);
        ticks--;
    end
endtask
```

Стоит отметить, что очень важным является то, что она ждёт не `posedge/negedge`, а именно такт, то есть данная функция не меняет чётности синхронизации. Это очень важно и много где используется. К сожалению, я не нашёл и не придумал способа не дублировать её код в каждом из модулей, так как если передавать `clk` по `input logic clk` аргумент, то он не будет обновляться внутри этой функции. Действительно в стандарте SystemVerilog на этот случай есть `ref logic clk`, что делало ровно то, что нам надо, но эта конструкция не поддерживается в Icarus Verilog 11, поэтому от её использования пришлось отказаться.

Про выравнивания структур

Стоит обратить внимание на ключевое слово `packed` в объявлениях структур. Выглядит это примерно так:

```
typedef struct packed {
    logic[cache_tag_size-1:0] tag;
    logic[cache_set_size-1:0] set;
    logic[cache_offset_size-1:0] offset;
} cacheAddr;
```

Такая конструкция гарантирует, что все данные структуры будут расположены последовательно, без каких-либо пропусков. Это позволяет значительно упростить себе присваивание значений для этой структуры, пользуясь конструкциями, описанными в следующем пункте.

Проблемы при записи в структуру

Из-за особенностей Verilog в массиве структур нельзя изменить отдельное поле, в этом случае при компиляции выдаётся следующая ошибка.

error: Array index expressions must be constant here.

Это решается тем, что записывать структуру целиком можно, это выглядит так:

```
lines[it] <= {1'b1, 1'b1, total_hits + total_misses,
curAddr.tag, buff_b};
```

В коде выше, конструкция с перечислением внутри фигурных скобок называется concatenation operator, и занимается тем, что конкатенирует всё внутри в один регистр **logic**. Интересным моментом является то, как написаны константы 1'b1, эта запись позволяет задаёт число в виде регистра заданной длины. В противном случае, запись 1 представляла бы собой регистра размера 32 и вся запись в структуру бы «поплыла».

Про блокирующие и неблокирующие присваивания

В коде выше используется нестандартное присваивание, которое называется неблокирующим. Оно применяется не сразу, а только в самом конце текущего такта, что позволяет делать следующие трюки. Вот код, который за такт меняет значение переменных a и b местами.

```
a <= b;
b <= a;
```

Его использование предохраняет от состояния «гонки» присвоения, так как гораздо проще спрогнозировать поведение несколько присвоений, которые происходят одновременно, чем нескольких последовательных. Помимо этого, код в котором меньше блокирующих присваиваний будет лучше преобразовываться в реальную схему.

В моей программе, я старался придерживаться следующей логики: переменные, которые используются для внутренней логики модуля и никак не взаимодействуют с другими модулями используют блокирующие присваивание, а значения на шинах и всё, что с ними связано — неблокирующее.

Протоколы взаимодействия модулей

Протокол работы с кэшем

При чтении из кэш у меня реализована следующая схема. На первом и втором такте передачи кэшу передаётся адрес, после чего владение шиной переходит кэшу. Спустя некоторое время, кэш на высокой синхронизации посылает значение **C1_RESPONSE**, одновременно с этим начиная передавать данные.

На схемах нарисован график синхронизации, на верхних значениях которого написаны действия кэша, а на нижних, соответственно, процессора в каждый такт синхронизации. Пунктирной линией обозначена смена владельца шины. В коде это реализовано таким образом, что процессор обрабатывает сразу после @ (**posedge** clk), а кэш наоборот, сразу после @ (**negedge** clk). Исключением из этого правила служит передача доступа от одной стороне другой, в этом случае сторона, которая прекращает владеть шиной, должна сообщить это не в свой такт, но в этот факт она просто снимает с себя владение шиной и ничего больше.

Приношу извинения за картинки от руки, inkscape неожиданно отказался работать должным образом.

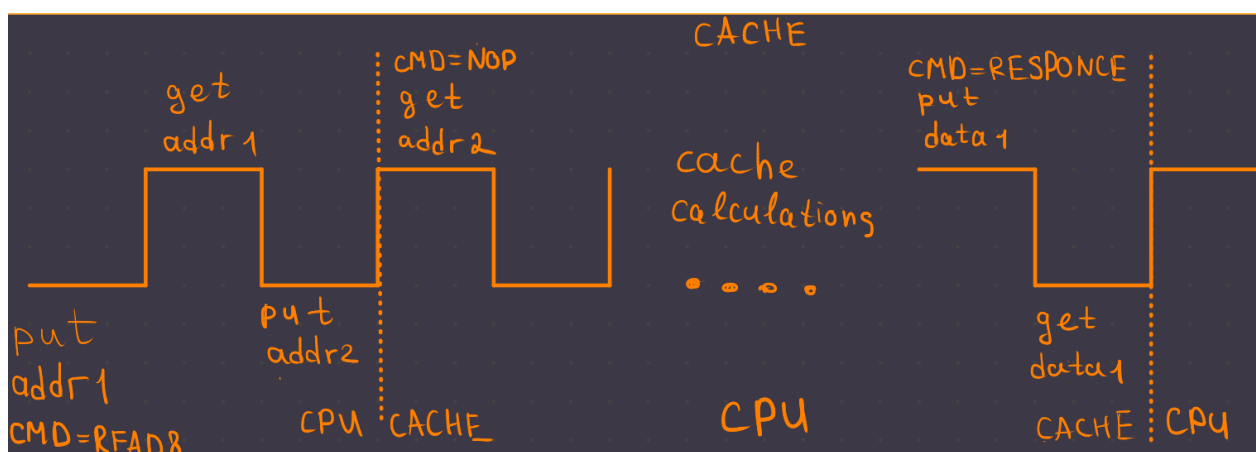


Рис. 1: C1_READ8

При запросе записи всё проходит почти аналогично. Обратите внимание, тут данные тоже передаются не сразу, а с задержкой в один такт от адреса. Позже оказалось, что это самое удачное решение, ведь мы тратим лишний такт, но решил всё таки оставить так.

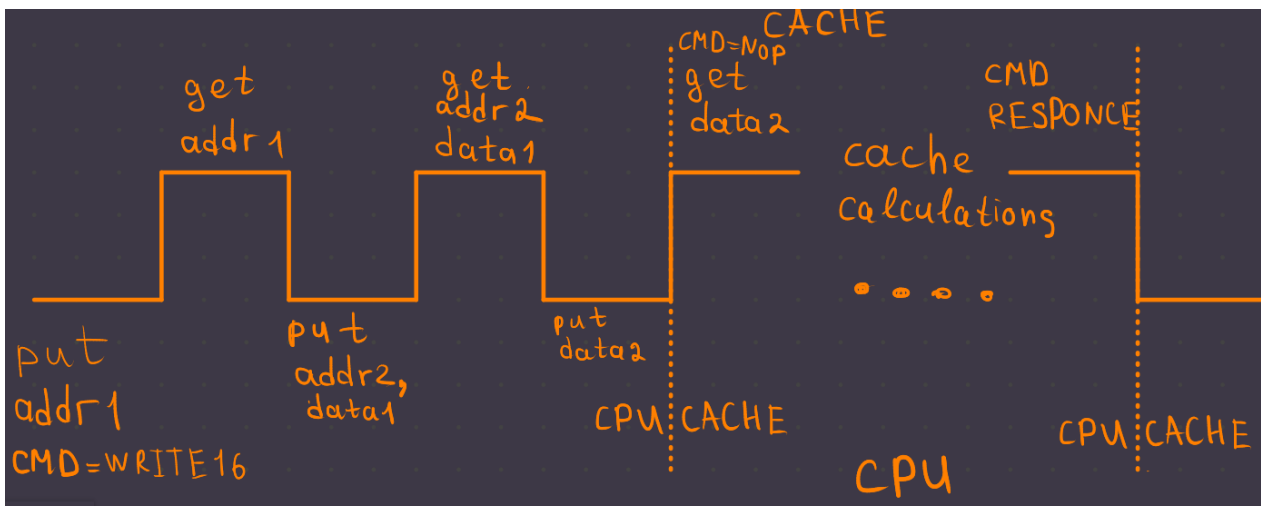


Рис. 2: C1_WRITE16

Протокол работы памятью

В ней всё очень похоже на работу кэша, но на всякий случай приведу схемы и для памяти. Сверху обозначены действия памяти, а снизу — кэша.

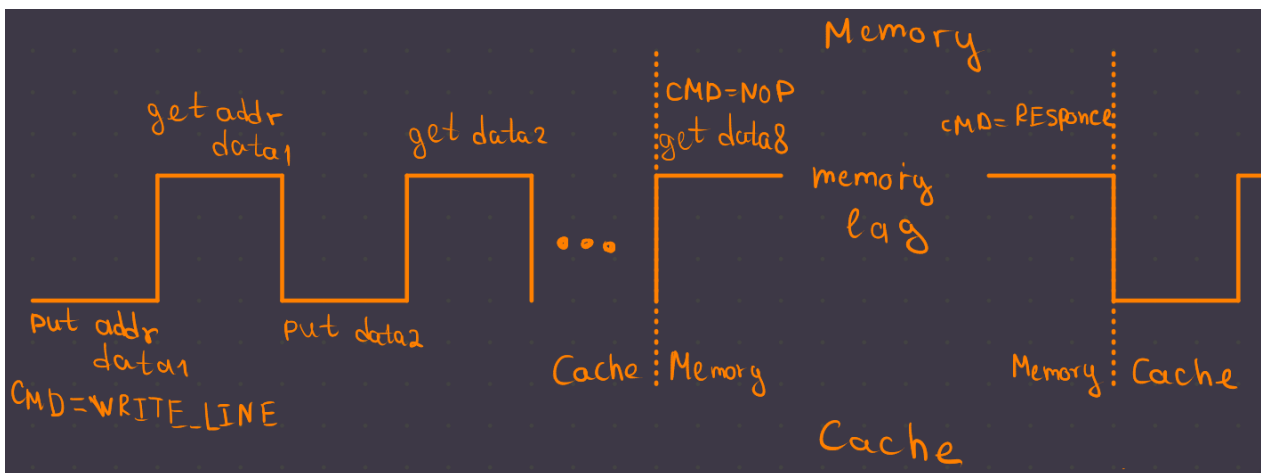


Рис. 3: C2_WRITELINE

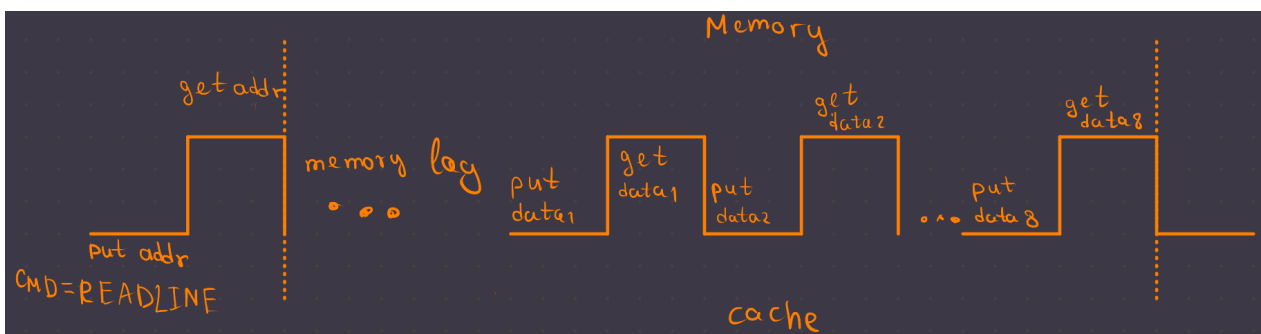


Рис. 4: C2_READLINE

Для наглядности также приведу код для считывания памяти с обеих сторон.

```

// Memory
always @(posedge clk) begin
    if (cmd_w == C2_READ_LINE) begin
        cmd <= C2_NOP;
        owner <= 1;
        skip(mem_feedback_time - 1);
        @(negedge clk)
        cmd <= C2_RESPONSE;
        for (it = addr * cache_line_size; it < addr *
cache_line_size + cache_line_size; it += data2_bus_size) begin
            for (byte_in_bus = 0; byte_in_bus < data2_bus_size;
byte_in_bus += 1) begin
                data[byte_in_bus * BITS_IN_BYTE +: BITS_IN_BYTE]
<= heap[it + byte_in_bus];
            end
            if (it + data2_bus_size >= addr * cache_line_size +
cache_line_size)
                @(posedge clk);
            else
                @(negedge clk);
        end
        owner <= 0;
    end
end

// Cache
task run_read(input logic[cache_set_size + cache_tag_size-1:0]
addr_, output logic[cache_line_size*BITS_IN_BYTE-1:0] data_,
output longint timing);
    @(negedge clk);
    owner <= 1;
    cmd <= C2_READ_LINE;
    addr <= addr_;
    timing = clk_time;
    @(posedge clk);
    owner <= 0;
    wait(cmd_w == C2_RESPONSE);
    @(posedge clk);
    timing = clk_time - timing;
    for (i = 0; i < cache_line_size; i += data2_bus_size) begin
        data_[i * BITS_IN_BYTE +: data2_bus_size * BITS_IN_BYTE]
<= data_w;
        if (i + data2_bus_size >= cache_line_size)
            @(negedge clk);
        else
            @(posedge clk);
    end
    owner <= 1;
    cmd <= C2_NOP;
endtask

```

Воспроизведение задачи на Verilog

В файле `testbench.sv` задано тестовое окружение для данной задачи, в котором выводится интересная статистика. В неё входит общее количество тактов работы задачи, количество тактов, затраченное на работу кэша, количество попаданий и промахов в кэш, а также время работы после каждой итерации внешнего цикла.

```
for (y = 0; y < M; y++) begin skip(); // loop
    skip(); // x init;
    ...
    ...
    $display("time: %d %t", y, timing);
    $fflush;
end

$display("Finish cpu run");
$display("Time: %t", timing);
$display("Cache time: %t", timing - skipped_time);
$display("Alu time: %t", skipped_time);
$display("Total hits: %d", total_hits);
$display("Total misses: %d", total_misses);
$finish;
```

Результат работы:

```
time: 0 83637
...
time: 62 5191123
time: 63 5272741
Finish cpu run
Time: 5272742
Cache time: 4274080
Alu time: 998662
Total hits: 228080
Total misses: 21520
```

Сравнение полученных результатов

Изначально прогнозированное время работы было меньше реально примерно на 20%. Это произошло главным образом из-за плохой реализации передачи данных по шине.

В коде на Verilog много где была сформирована модель, что принимающая сторона синхронизируется по **posedge** `clk`, а отдающая по **negedge** `clk` или наоборот. Из-за этого очередная операция считывания могла быть начата на такт позже из-за того, что на данном такте **negedge** `clk`, а эта схема синхронизируется по **posedge** `clk` перед отправкой.

Но описанная выше проблема была исправлена после исправления не оптимальных мест и изменения логики синхронизации. В этом мне очень помогли написанные

тесты. После этих изменений, время, полученное во время симуляции кода на Verilog в точности совпало с прогнозом.

Листининг кода

Код на C++ для оценки времени работы

```
1  #include <iostream>
2  #include <vector>
3  #include <algorithm>
4  #include <numeric>
5  #include <bitset>
6  #include <cassert>
7  #include <fstream>
8  #include <tuple>
9
10 using namespace std;
11
12 typedef unsigned long long Addr;
13
14 enum Word { WORD=8, DWORD=16, QWORD=32 };
15
16 struct CacheLine {
17     const int size;
18     bool valid = 0;
19     bool dirty = 0;
20     Addr tag = -1;
21     int last_call = 0;
22     CacheLine& operator = (CacheLine oth) {
23         assert(size == oth.size);
24         tie(valid, dirty, tag, last_call) = tie(oth.valid, oth.dirty, oth.tag, oth.last_call);
25         return *this;
26     }
27 };
28
29 struct CacheSet {
30 private:
31     vector<CacheLine> lines;
32     const int line_size;
33     int total_mem_pushes = 0;
34 public:
35     CacheSet(int way, int line_size) : line_size(line_size) {
36         lines = vector<CacheLine>(way, CacheLine{line_size});
37     }
38
39     bool exists(Addr tag) {
40         return find_by_tag(tag) != lines.end();
41     }
42
43     void read(Addr tag, int call) {
44         auto it = find_by_tag(tag);
45         if (it != lines.end()) {
46             it->last_call = call;
47         } else {
48             *(find_LRO()) = CacheLine {line_size, 1, 0, tag, call};
49         }
50     }
51
52     void write(Addr tag, int call) {
53         auto it = find_by_tag(tag);
54         if (it != lines.end()) {
55             it->last_call = call;
56             it->dirty = 1;
57         } else {
58             total_mem_pushes += find_LRO()->dirty;
59             *(find_LRO()) = CacheLine {line_size, 1, 1, tag, call};
60         }
61     }
62
63     int getMemPushes() const {
64         return total_mem_pushes;
65     }
66 }
```

```

67 private:
68     vector<CacheLine>::iterator find_by_tag(Addr tag) {
69         return find_if(lines.begin(), lines.end(), [&] (const CacheLine& line) {
70             return line.valid && line.tag == tag;
71         });
72     };
73 }
74
75 vector<CacheLine>::iterator find_LRO() {
76     return min_element(lines.begin(), lines.end(), [] (const CacheLine& a, const CacheLine& b) {
77         return tie(a.valid, a.last_call) < tie(b.valid, b.last_call);
78     });
79 }
80 };
81
82
83 struct Cache {
84 private:
85     const int sets_cnt;
86     const int way;
87     const int line_size;           // byte
88     const int data_bus_size;      // bits
89     const int mem_size;           // byte
90
91     int total_hits = 0;
92     int total_misses = 0;
93     int total_time = 0;
94     int calls = 0;
95
96     vector<CacheSet> sets;
97
98     struct InnerAddr {
99         Addr tag, set, offset;
100     };
101 public:
102     Cache(int sets_cnt, int way, int line_size, int data_bus_size, int mem_size) :
103         sets_cnt(sets_cnt), way(way), line_size(line_size), data_bus_size(data_bus_size), mem_size(mem_size)
104     {
105         sets = vector<CacheSet>(sets_cnt, CacheSet(way, line_size));
106     }
107
108     void read(Word word, Addr addr) {
109         calls++;
110         auto [tag, set, offset] = split_addr(addr);
111         if (sets[set].exists(tag)) {
112             total_hits++;
113             total_time += 6; // Cache lag
114             total_time += transfer_lag(word) + 1; // Cache -> Cpu
115             // 8-9
116         } else {
117             total_misses++;
118             total_time += 4 + 100; // Cache + Mem lag
119             total_time += transfer_lag(WORD * line_size) + 1; // Mem -> Cache
120             total_time += transfer_lag(word); // Cache -> Cpu
121             // 114-115
122         }
123         sets[set].read(tag, calls);
124     };
125
126     void write(Word word, Addr addr) {
127         calls++;
128         auto [tag, set, offset] = split_addr(addr);
129         if (sets[set].exists(tag)) {
130             total_hits++;
131             total_time += 6; // Cache lag
132             total_time += 1; // Cache->Cpu response
133         } else {
134             total_misses++;
135             total_time += 4 + 100; // Cache + Mem lag
136             total_time += transfer_lag(WORD * line_size); // Mem -> Cache
137             total_time += 1; // Cache->Cpu response
138         }
139         sets[set].write(tag, calls);
140     };
141
142     int get_hits() const {
143         return total_hits;
144     }
145
146     int get_misses() const {

```

```

146     return total_misses;
147 }
148
149 int get_time() const {
150     int total_pushes = accumulate(sets.begin(), sets.end(), 0,
151     [] (int acc, const CacheSet& el) {
152         return acc + el.getMemPushes();
153     });
154 };
155     return total_time + total_pushes * (100 + 1);
156 }
157
158 private:
159     InnerAddr split_addr(Addr addr) const {
160         return {
161             (addr & ((mem_size - 1) & ~(sets_cnt * line_size - 1))) / (sets_cnt * line_size),
162             (addr & (sets_cnt - 1) * line_size) / line_size,
163             addr & (line_size - 1)
164         };
165     }
166
167     int transfer_lag(int data) const {
168         return (data + data_bus_size - 1) / data_bus_size;
169     }
170 };
171
172 struct PseudoAllocator {
173     Addr last_allocated_addr = 0;
174
175     Addr allocate(int n, int m, int elem_sz) {
176         int res = last_allocated_addr;
177         last_allocated_addr += n * m * elem_sz;
178         return res;
179     }
180 };
181
182 int main() {
183     // 512Kb == [0x00000...0x7ffff]
184     Cache cache(32, 2, 16, 16, 512 * 1024);
185     PseudoAllocator alloc;
186
187     const int M = 64;
188     const int N = 60;
189     const int K = 32;
190
191     int8_t* const a = (int8_t*) alloc.allocate(M, K, sizeof(int8_t)); // a[M][K];
192     int16_t* const b = (int16_t*) alloc.allocate(K, N, sizeof(int16_t)); // b[K][N];
193     int32_t* const c = (int32_t*) alloc.allocate(M, N, sizeof(int32_t)); // c[M][N];
194
195     int8_t* pa = a;
196     int32_t* pc = c;
197
198     for (int y = 0; y < M; y++) {
199         for (int x = 0; x < N; x++) {
200             int16_t* pb = b;
201             for (int k = 0; k < K; k++) {
202                 cache.read(WORD, (Addr) (pa + k));
203                 cache.read(DWORD, (Addr) (pb + x));
204                 pb += N;
205             }
206             cache.write(QWORD, (Addr) (pc + x));
207         }
208         pa += K;
209         pc += N;
210     }
211
212     cout << "HITS: " << cache.get_hits() << "\nMISSES: " << cache.get_misses() << "\nTOTAL TIME: " <<
213     cache.get_time() << endl;
214     cout << "RATE: " << 1.0 * (cache.get_hits()) / (cache.get_hits() + cache.get_misses()) << endl;
215
216     return 0;
217 }

```

Listing 1: main.cpp

Код на Verilog для симуляции схемы

```
1  `ifndef PARAMETERS_GUARD
2  `define PARAMETERS_GUARD
3
4  parameter BITS_IN_BYTE = 8;
5  parameter INF = 1000000000;
6
7  parameter cache_line_size = 16;
8  parameter cache_way = 2;
9  parameter cache_tag_size = 10;
10 parameter cache_set_size = 5;
11 parameter cache_offset_size = 4;
12
13 parameter addr1_bus_size = 2;
14 parameter data1_bus_size = 2;
15
16 parameter data2_bus_size = 2;
17 parameter addr2_bus_size = 2;
18
19 parameter mem_feedback_time = 100;
20 parameter cache_hit_time = 6;
21 parameter cache_miss_time = 4 + mem_feedback_time + cache_line_size / data2_bus_size;
22
23 parameter _SEED = 225526;
24
25 typedef enum logic[1:0] {
26     C2_NOP=2'b00,
27     C2_RESPONSE=2'b01,
28     C2_READ_LINE=2'b10,
29     C2_WRITE_LINE=2'b11
30 } mem_command_codes;
31
32 typedef enum logic[2:0] {
33     C1_NOP=3'b000,
34     C1_READ8=3'b001,
35     C1_READ16=3'b010,
36     C1_READ32=3'b011,
37     C1_INVALIDATE_LINE=3'b100,
38     C1_WRITE8=3'b101,
39     C1_WRITE16=3'b110,
40     C1_WRITE32=3'b111
41 } cpu_cache_command_codes;
42
43 typedef enum logic[2:0] {
44     C1_RESPONSE=3'b111
45 } cache_cpu_command_codes;
46 `endif
```

Listing 2: parameters.sv

```
1  `ifndef CLOCK_GUARD
2  `define CLOCK_GUARD
3
4  module Clock(output logic clk, output longint timing);
5
6  always #1 clk = ~clk;
7
8  always #2 timing++;
9
10 initial begin
11     timing = 0;
12     clk = 0;
13 end
14
15 endmodule
16
17 `endif
```

Listing 3: clock.sv

```

1  `ifndef MEMORY_GUARD
2  `define MEMORY_GUARD
3
4  `include "parameters.sv"
5  `include "clock.sv"
6
7  module Memory
8      #(
9          parameter mem_size = 512 * 1024
10     )
11     (
12         input wire clk, reset, dump,
13         input wire[addr2_bus_size*BITS_IN_BYTE-1:0] addr,
14         inout wire[data2_bus_size*BITS_IN_BYTE-1:0] data_w,
15         inout wire[1:0] cmd_w
16     );
17
18     logic[BITS_IN_BYTE-1:0] heap[0:mem_size-1];
19     integer SEED = _SEED;
20     integer i = 0;
21
22     bit owner = 0;
23     logic[1:0] cmd;
24     logic[data2_bus_size*BITS_IN_BYTE-1:0] data;
25
26     assign cmd_w = owner ? cmd : 2'bzz;
27     assign data_w = owner ? data : {data2_bus_size*BITS_IN_BYTE{1'bz}};
28
29     always @(posedge reset) begin
30         for (i = 0; i < mem_size; i += 1) begin
31             heap[i] <= $random(SEED)>>16;
32         end
33     end
34
35     integer fd;
36     always @(posedge dump) begin
37         fd = $fopen("mem.dump", "w");
38         for (i = 0; i < mem_size; i += 1) begin
39             $display(fd, "[%d] %d", i, heap[i]);
40         end
41         $display("Memory has dumped to mem.dump");
42         $fclose(fd);
43     end
44
45     task automatic skip(input longint ticks = 1);
46         logic enter_clk = clk;
47         while (ticks > 0) begin
48             wait(clk != enter_clk);
49             wait(clk == enter_clk);
50             ticks--;
51         end
52     endtask
53
54     integer it, byte_in_bus;
55
56     always @(posedge clk) begin
57         if (cmd_w == C2_WRITE_LINE) begin
58             for (it = addr * cache_line_size; it < addr * cache_line_size + cache_line_size; it +=
59                 data2_bus_size) begin
60                 for (byte_in_bus = 0; byte_in_bus < data2_bus_size; byte_in_bus += 1) begin
61                     heap[it + byte_in_bus] <= data_w[byte_in_bus * BITS_IN_BYTE +: BITS_IN_BYTE];
62                 end
63                 @(posedge clk);
64             end
65             skip(mem_feedback_time - cache_line_size / data2_bus_size - 1);
66             @(negedge clk);
67             cmd <= C2_RESPONSE;
68             owner <= 1;
69             @(posedge clk);
70             owner <= 0;
71         end
72     end
73
74     always @(posedge clk) begin
75         if (cmd_w == C2_READ_LINE) begin
76             cmd <= C2_NOP;
77             owner <= 1;
78             skip(mem_feedback_time - 1);
79             @(negedge clk)
80             cmd <= C2_RESPONSE;

```

```

80         for (it = addr * cache_line_size; it < addr * cache_line_size + cache_line_size; it +=
81             data2_bus_size) begin
82             for (byte_in_bus = 0; byte_in_bus < data2_bus_size; byte_in_bus += 1) begin
83                 data[byte_in_bus * BITS_IN_BYTE +: BITS_IN_BYTE] <= heap[it + byte_in_bus];
84             end
85             if (it + data2_bus_size >= addr * cache_line_size + cache_line_size)
86                 @(posedge clk);
87             else
88                 @(negedge clk);
89             end
90             owner <= 0;
91         end
92     endmodule
93
94 module MemoryDriver
95 (
96     input logic clk,
97     input longint clk_time,
98     input logic reset,
99     input logic m_dump,
100     output logic[addr2_bus_size*BITS_IN_BYTE-1:0] addr_w,
101     inout logic[data2_bus_size*BITS_IN_BYTE-1:0] data_w,
102     inout logic[1:0] cmd_w
103 );
104
105 logic[1:0] cmd = C2_NOP;
106 logic[addr2_bus_size*BITS_IN_BYTE-1:0] addr;
107 logic[data2_bus_size*BITS_IN_BYTE-1:0] data;
108 bit owner = 1;
109
110 assign cmd_w = owner ? cmd : 2'bzz;
111 assign addr_w = addr;
112 assign data_w = owner ? data : {data2_bus_size*BITS_IN_BYTE{1'bz}};
113
114 integer i;
115
116 task run_read(input logic[cache_set_size + cache_tag_size-1:0] addr_, output
117 logic[cache_line_size*BITS_IN_BYTE-1:0] data_, output longint timing);
118     @(negedge clk);
119     owner <= 1;
120     cmd <= C2_READ_LINE;
121     addr <= addr_;
122     timing = clk_time;
123     @(posedge clk);
124     owner <= 0;
125     wait(cmd_w == C2_RESPONSE);
126     @(posedge clk);
127     timing = clk_time - timing;
128     for (i = 0; i < cache_line_size; i += data2_bus_size) begin
129         data[i * BITS_IN_BYTE +: data2_bus_size * BITS_IN_BYTE] <= data_w;
130         if (i + data2_bus_size >= cache_line_size)
131             @(negedge clk);
132         else
133             @(posedge clk);
134     end
135     owner <= 1;
136     cmd <= C2_NOP;
137 endtask
138
139 task run_write(input int addr_, input logic[cache_line_size*BITS_IN_BYTE-1:0] data_, output longint
140 timing);
141     @(negedge clk);
142     timing = clk_time;
143     owner <= 1;
144     cmd <= C2_WRITE_LINE;
145     addr <= addr_;
146     for (i = 0; i < cache_line_size / data2_bus_size; i += 1) begin
147         data <= data_[i * data2_bus_size * BITS_IN_BYTE +: data2_bus_size * BITS_IN_BYTE];
148         @(negedge clk);
149     end
150     owner <= 0;
151     wait(cmd_w == C2_RESPONSE);
152     timing = clk_time - timing;
153     cmd <= C2_NOP;
154     owner <= 1;
155 endtask
156 endmodule
`endif

```

Listing 4: mem.sv

```

1  `ifndef CACHE_GUARD
2  `define CACHE_GUARD
3
4  `include "parameters.sv"
5  `include "mem.sv"
6  `include "clock.sv"
7
8  module Cache (
9      input wire                                clk, reset, dump,
10     input wire[addr1_bus_size*BITS_IN_BYTE-1:0] addr_cpu_w,
11     inout wire[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu_w,
12     inout wire[2:0] cmd_cpu_w,
13     output wire[addr2_bus_size*BITS_IN_BYTE-1:0] addr_mem_w,
14     inout wire[data2_bus_size*BITS_IN_BYTE-1:0] data_mem_w,
15     inout wire[1:0] cmd_mem_w,
16     output integer total_hits = 0,
17     output integer total_misses = 0
18 );
19
20 localparam cache_sets_count = 1 << cache_set_size;
21
22 bit owner_cpu = 0;
23 logic[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu;
24 logic[2:0] cmd_cpu;
25 assign data_cpu_w = owner_cpu ? data_cpu : {data1_bus_size*BITS_IN_BYTE{1'bz}};
26 assign cmd_cpu_w = owner_cpu ? cmd_cpu : 3'bzzz;
27
28 MemoryDriver mem_driver(clk, 64'b0, reset, 1'b0, addr_mem_w, data_mem_w, cmd_mem_w);
29
30 task run_mem_write(input logic[cache_set_size + cache_tag_size-1:0] addr_, input
31     ↪ logic[cache_line_size*BITS_IN_BYTE-1:0] data_);
32     longint timing; // Unusable blackhole
33     mem_driver.run_write(addr_, data_, timing);
34 endtask
35
36 task run_mem_read(input logic[cache_set_size + cache_tag_size-1:0] addr_, output
37     ↪ logic[cache_line_size*BITS_IN_BYTE-1:0] data_);
38     longint timing; // Unusable blackhole
39     mem_driver.run_read(addr_, data_, timing);
40 endtask
41
42 typedef struct packed {
43     logic valid;
44     logic dirty;
45     integer last_update;
46     logic[cache_tag_size-1:0] tag;
47     logic[cache_line_size*BITS_IN_BYTE-1:0] data;
48 } cacheLine;
49
50 typedef struct packed {
51     logic[cache_tag_size-1:0] tag;
52     logic[cache_set_size-1:0] set;
53     logic[cache_offset_size-1:0] offset;
54 } cacheAddr;
55
56 cacheLine[cache_way * cache_sets_count-1:0] lines;
57
58 function int bytes_cnt_from_cmd(logic[2:0] cmd);
59     case (cmd & 3'b011)
60         3'b001 : return 1;
61         3'b010 : return 2;
62         3'b011 : return 4;
63     endcase
64 endfunction
65
66 task automatic skip(input longint ticks = 1);
67     logic enter_clk = clk;
68     while (ticks > 0) begin
69         wait(clk != enter_clk);
70         wait(clk == enter_clk);
71         ticks--;
72     end
73 endtask
74
75 always @(posedge reset) begin

```

```

74     for (i = 0; i < cache_way * cache_sets_count; i += 1) begin
75         lines[i] <= 0;
76     end
77 end
78
79 integer action_word, it, i, j;
80 logic[cache_line_size * BITS_IN_BYTE-1:0] buff_a, buff_b;
81 cacheAddr curAddr;
82
83 always @(posedge clk) begin
84     if (owner_cpu == 0 && (cmd_cpu_w == C1_WRITE8 || cmd_cpu_w == C1_WRITE16 || cmd_cpu_w ==
85 C1_WRITE32)) begin
86         // READ 1-ST
87         action_word <= bytes_cnt_from_cmd(cmd_cpu_w);
88         {curAddr.tag, curAddr.set} <= addr_cpu_w;
89         // READ 2-ND
90         @(posedge clk);
91         curAddr.offset = addr_cpu_w;
92         // Read data
93         for (j = 0; j < action_word; j += data1_bus_size) begin
94             for (i = 0; i + j < action_word && i < data1_bus_size; i++) begin
95                 buff_a[(curAddr.offset + j + i) * BITS_IN_BYTE +: BITS_IN_BYTE] <= data_cpu_w[i *
96 BITS_IN_BYTE +: BITS_IN_BYTE];
97             end
98             if (j + data1_bus_size >= action_word)
99                 @(negedge clk);
100             else
101                 @(posedge clk);
102         end
103         owner_cpu <= 1;
104         cmd_cpu <= C1_NOP;
105
106         // Searching line
107         it = search_by_addr_or_empty(curAddr, lines);
108         // Uploading from memory if needed
109         if (it == -1 || !lines[it].valid || lines[it].tag != curAddr.tag) begin
110             total_misses++;
111             skip(action_word != 4);
112             run_mem_read({curAddr.tag, curAddr.set}, buff_b);
113         end else begin
114             total_hits++;
115             skip(3 + (action_word != 4));
116             buff_b <= lines[it];
117         end
118
119         // Merging buffers
120         for (i = 0; i < action_word; i++) begin
121             buff_b[(curAddr.offset + i) * BITS_IN_BYTE +: BITS_IN_BYTE] = buff_a[(curAddr.offset + i)
122 * BITS_IN_BYTE +: BITS_IN_BYTE];
123         end
124
125         // Purge if needed
126         if (it == -1) begin
127             it = find_lru(curAddr.set, lines);
128             if (lines[it].dirty) begin
129                 $display("PURGE");
130                 run_mem_write({lines[it].tag, curAddr.set}, lines[it]);
131             end
132         end
133
134         lines[it] <= {1'b1, 1'b1, total_hits + total_misses, curAddr.tag, buff_b};
135
136         cmd_cpu <= C1_RESPONSE;
137         @(posedge clk);
138         owner_cpu <= 0;
139     end
140 end
141
142 function int search_by_addr_or_empty(
143     input cacheAddr addr,
144     input cacheLine[cache_way * cache_sets_count-1:0] lines
145 );
146 it = -1;
147 for (i = addr.set * cache_way; i < addr.set * cache_way + cache_way; ++i) begin
148     if (lines[i].valid && lines[i].tag == addr.tag)
149         it = i;
150     if (!lines[i].valid && it == -1)
151         it = i;
152 end
153 return it;

```



```

151     endfunction
152
153     function int find_lru(input logic[cache_set_size-1:0] set, input cacheLine[cache_way *
154     cache_sets_count-1:0] lines);
155         integer it_time;
156         it_time = INF;
157         it = -1;
158         for (i = set * cache_way; i < set * cache_way + cache_way; ++i) begin
159             if (it_time > lines[i].last_update) begin
160                 it = i;
161                 it_time = lines[i].last_update;
162             end
163         end
164         return it;
165     endfunction
166
167     logic[cache_line_size * BITS_IN_BYTE-1:0] buff;
168     always @(posedge clk) begin
169         if (owner_cpu == 0 && (cmd_cpu_w == C1_READ8 || cmd_cpu_w == C1_READ16 || cmd_cpu_w ==
170         C1_READ32)) begin
171             // READ 1-ST
172             action_word <= bytes_cnt_from_cmd(cmd_cpu_w);
173             {curAddr.tag, curAddr.set} <= addr_cpu_w;
174             // READ 2-ND
175             @(posedge clk);
176             curAddr.offset = addr_cpu_w;
177             owner_cpu <= 1;
178             cmd_cpu <= C1_NOP;
179
180             it = search_by_addr_or_empty(curAddr, lines);
181
182             if (it == -1) begin
183                 it = find_lru(curAddr.set, lines);
184             end
185             if (!lines[it].valid || lines[it].tag != curAddr.tag) begin
186                 total_misses++;
187                 skip(2);
188                 run_mem_read({curAddr.tag, curAddr.set}, buff);
189                 lines[it] <= {1'b1, 1'b0, total_hits + total_misses, curAddr.tag, buff};
190             end else begin
191                 total_hits++;
192                 skip(5);
193                 buff = lines[it].data;
194                 lines[it] <= {lines[it].valid, lines[it].dirty, total_hits + total_misses, lines[it].tag,
195                 buff};
196             end
197
198             cmd_cpu <= C1_RESPONSE;
199             for (j = 0; j < action_word; j += data1_bus_size) begin
200                 for (i = 0; i + j < action_word && i < data1_bus_size; i++) begin
201                     data_cpu[i * BITS_IN_BYTE +: BITS_IN_BYTE] = buff[(curAddr.offset + j + i) *
202                     BITS_IN_BYTE +: BITS_IN_BYTE];
203                 end
204                 if (j + data1_bus_size >= action_word)
205                     @(posedge clk); // last iteration
206                 else
207                     @(negedge clk);
208             end
209             owner_cpu <= 0;
210         end
211     end
212 endmodule
213
214 module CacheDriver
215 (
216     input logic clk,
217     input longint clk_time,
218     input logic reset,
219     input logic c_dump,
220     output logic[addr1_bus_size*BITS_IN_BYTE-1:0] addr_cpu_w,
221     inout logic[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu_w,
222     inout logic[2:0] cmd_cpu_w
223 );
224
225 bit owner_cpu = 1;
226 logic[2:0] cmd_cpu = C1_NOP;
227 logic[addr1_bus_size*BITS_IN_BYTE-1:0] addr_cpu;
228 logic[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu;
229 assign addr_cpu_w = addr_cpu;

```

```

227 assign data_cpu_w = owner_cpu ? data_cpu : {data1_bus_size*BITS_IN_BYTE{1'bz}};
228 assign cmd_cpu_w = owner_cpu ? cmd_cpu : 3'bzzz;
229
230 task run_read(
231     input logic[cache_tag_size + cache_offset_size + cache_set_size - 1 : 0] addr,
232     input logic[2:0] cmd,
233     output logic[BITS_IN_BYTE*cache_line_size-1:0] data,
234     output longint timing
235 );
236     logic[BITS_IN_BYTE*data1_bus_size-1:0] local_buff;
237     @(negedge clk);
238     timing = clk_time;
239     cmd_cpu <= cmd;
240     addr_cpu <= addr[cache_offset_size +: cache_set_size + cache_tag_size];
241     @(negedge clk);
242     addr_cpu <= addr[0 +: cache_offset_size];
243     @(posedge clk);
244     owner_cpu <= 0;
245     wait(cmd_cpu_w == C1_RESPONSE); // response on posedge!
246     /* @(posedge clk); */
247     timing = clk_time - timing;
248     case (cmd)
249         C1_READ8, C1_READ16 : begin
250             data = data_cpu_w;
251         end
252         C1_READ32 : begin
253             local_buff <= data_cpu_w;
254             @(posedge clk);
255             data = {data_cpu_w[0 +: data1_bus_size*BITS_IN_BYTE], local_buff[0 +:
256                 data1_bus_size*BITS_IN_BYTE]};
257             end
258             default : begin
259                 $display("Incorrect run_read cmd: %d", cmd);
260                 $finish;
261             end
262         endcase
263         owner_cpu <= 1;
264         cmd_cpu <= C1_NOP;
265         @(negedge clk);
266     endtask
267
268 task run_write(
269     input logic[cache_tag_size + cache_offset_size + cache_set_size - 1 : 0] addr,
270     input logic[2:0] cmd,
271     input logic[BITS_IN_BYTE*cache_line_size - 1:0] data,
272     output longint timing
273 );
274     /* $monitor("time %t, owner: %b, clk: %b cmd_w: %b addr_w: %b data_w: %b", clk_time, owner_cpu,
275     clk, cmd_cpu_w, addr_cpu_w, data_cpu_w); */
276     @(negedge clk);
277     timing = clk_time;
278     cmd_cpu <= cmd;
279     addr_cpu <= addr[cache_offset_size +: cache_set_size + cache_tag_size];
280     @(negedge clk);
281     addr_cpu <= addr[0 +: cache_offset_size];
282     case (cmd)
283         C1_WRITE8, C1_WRITE16 : begin
284             data_cpu <= data;
285         end
286         C1_WRITE32 : begin
287             data_cpu <= data[0 +: data1_bus_size * BITS_IN_BYTE];
288             @(negedge clk);
289             data_cpu <= data[data1_bus_size*BITS_IN_BYTE +: data1_bus_size*BITS_IN_BYTE];
290         end
291         default : begin
292             $display("Incorrect run_write cmd: %d", cmd);
293             $finish;
294         end
295     endcase
296     owner_cpu <= 0;
297     @(negedge clk); // ???
298     wait(cmd_cpu_w == C1_RESPONSE);
299     timing = clk_time - timing;
300     owner_cpu <= 1;
301     cmd_cpu <= C1_NOP;
302 endtask
303 endmodule
304 `endif

```

Listing 5: cache.sv

Код тестирования в среде Verilog

```
1 `include "cache.sv"
2 `include "clock.sv"
3
4 module Testbench;
5
6     integer skipped_time = 0;
7     task automatic skip(input longint ticks = 1);
8         logic enter_clk = clk;
9         skipped_time += ticks;
10        while (ticks > 0) begin
11            wait(clk != enter_clk);
12            wait(clk == enter_clk);
13            ticks--;
14        end
15    endtask
16
17    integer total_hits, total_misses;
18    longint timing;
19    logic clk, reset = 0, c_dump = 0, m_dump = 0;
20
21    wire[addr2_bus_size*BITS_IN_BYTE-1:0] addr_mem_w;
22    wire[data2_bus_size*BITS_IN_BYTE-1:0] data_mem_w;
23    wire[1:0] cmd_mem_w;
24    wire[addr1_bus_size*BITS_IN_BYTE-1:0] addr_cpu_w;
25    wire[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu_w;
26    wire[2:0] cmd_cpu_w;
27
28    Clock cloker(clk, timing);
29    Memory mem(clk, reset, m_dump, addr_mem_w, data_mem_w, cmd_mem_w);
30    Cache cache(clk, reset, c_dump, addr_cpu_w, data_cpu_w, cmd_cpu_w, addr_mem_w, data_mem_w,
31    cmd_mem_w, total_hits, total_misses);
32    CacheDriver cpu(clk, timing, reset, c_dump, addr_cpu_w, data_cpu_w, cmd_cpu_w);
33
34    logic[BITS_IN_BYTE-1:0] wbuff;
35    logic[2 * BITS_IN_BYTE-1:0] dbuff;
36    logic[4 * BITS_IN_BYTE-1:0] qbuff;
37
38    localparam M = 64;
39    localparam N = 60;
40    localparam K = 32;
41    integer a_addr = 0;
42    integer b_addr = a_addr + M * K;
43    integer c_addr = b_addr + K * N * 2;
44
45    integer y, x, k;
46    integer pa, pb, pc;
47    integer s;
48    integer blackhole;
49
50    initial begin
51        reset <= 1;
52        skip();
53        reset <= 0;
54        skip();
55
56        pa = a_addr; skip(); // init pa
57        pc = c_addr; skip(); // init pc
58
59        skip(); // y init;
60        for (y = 0; y < M; y++) begin skip(); // loop
61            skip(); // x init;
62            for (x = 0; x < N; x++) begin skip(); // loop
63                pb = b_addr; skip(); // init pb
64                s = 0; skip(); // init s
65                skip(); // k init;
66                for (k = 0; k < K; k++) begin skip(); // loop
67                    // s += pa[k] * pb[x] begin
68                    cpu.run_read(pa + k, C1_READ8, wbuff, blackhole);
69                    cpu.run_read(pb + x * 2, C1_READ16, dbuff, blackhole);
```

```

69         qbuff = wbuff * dbuff; skip(5); // (*)
70         s = s + qbuff; skip(); // (+)
71         // s += pa[k] * pb[x] end
72         pb = pb + N * 2; skip(); // (+)
73     end
74     cpu.run_write(pc + x * 4, C1_WRITE32, s, blackhole);
75 end
76 pa = pa + K; skip(); // (+)
77 pc = pc + N * 4; skip(); // (+)
78 $display("time: %d %t", y, timing);
79 $fflush;
80 end
81 skip(); // function exit
82
83 $display("Finish cpu run");
84 $display("Time: %t", timing);
85 $display("Cache time: %t", timing - skipped_time);
86 $display("Alu time: %t", skipped_time);
87 $display("Total hits: %d", total_hits);
88 $display("Total misses: %d", total_misses);
89 $finish;
90
91 /* Finish cpu run */
92 /* Time:          5272742 */
93 /* Cache time     4274080 */
94 /* Alu time       998662 */
95 /* Total hits:    228080 */
96 /* Total misses:  21520 */
97
98 end
99 endmodule

```

Listing 6: testbench.sv

```

1  `include "mem.sv"
2
3  module MemoryTestbench;
4
5      logic reset = 0, m_dump = 0, clk;
6      wire[addr2_bus_size*BITS_IN_BYTE-1:0] addr_w;
7      wire[data2_bus_size*BITS_IN_BYTE-1:0] data_w;
8      wire[1:0] cmd_w;
9      longint clk_time;
10
11      Clock clock(clk, clk_time);
12      Memory #(64) mem(clk, reset, m_dump, addr_w, data_w, cmd_w);
13      MemoryDriver driver(clk, clk_time, reset, m_dump, addr_w, data_w, cmd_w);
14
15      integer i;
16
17      logic[cache_line_size * BITS_IN_BYTE - 1 : 0] buff_a, buff_b;
18      longint timing = 100;
19      integer test_addr = 5;
20      logic[cache_line_size * BITS_IN_BYTE - 1 : 0] test_payload = {cache_line_size{$random()}};
21
22      always @(timing)
23          if (timing != mem_feedback_time)
24              $display("Memory timing error! Expected: %d, Real %d", mem_feedback_time, timing);
25
26      initial begin
27          reset <= 1;
28          #1;
29          reset <= 0;
30          #1;
31
32          begin : TEST_ZERO_ADDR
33              driver.run_read(0, buff_a, timing);
34              driver.run_write(0, test_payload, timing);
35              driver.run_read(0, buff_b, timing);
36              if (buff_b != test_payload)
37                  $display("Memory correctness unit test failed, real: %b expected: %b", buff_b,
38 test_payload);
39          end
40
41          begin : CHECK_CUSTOM_ADDR
42              driver.run_write(test_addr, buff_a, timing);
43              driver.run_read(test_addr, buff_b, timing);
44              if (buff_a != buff_b)

```

```

44         $display("Memory correctness unit test failed, real: %b expected: %b", buff_b, buff_a);
45     end
46
47     $display("Finish memory testing");
48     $finish;
49 end
50 endmodule

```

Listing 7: mem_testbench.sv

```

1  `include "cache.sv"
2
3  module CacheTestbench;
4
5      logic                                clk, reset=0, c_dump=0, m_dump=0;
6      wire[addr1_bus_size*BITS_IN_BYTE-1:0] addr_cpu_w;
7      wire[data1_bus_size*BITS_IN_BYTE-1:0] data_cpu_w;
8      wire[2:0]                             cmd_cpu_w ;
9      wire[addr2_bus_size*BITS_IN_BYTE-1:0] addr_mem_w;
10     wire[data2_bus_size*BITS_IN_BYTE-1:0] data_mem_w;
11     wire[1:0]                             cmd_mem_w ;
12     integer total_hits, total_misses;
13
14     longint clk_time;
15     Clock clker(clk, clk_time);
16     Memory mem(clk, reset, m_dump, addr_mem_w, data_mem_w, cmd_mem_w);
17     Cache cache(clk, reset, c_dump, addr_cpu_w, data_cpu_w, cmd_cpu_w, addr_mem_w, data_mem_w,
18     cmd_mem_w, total_hits, total_misses);
19     CacheDriver driver(clk, clk_time, reset, c_dump, addr_cpu_w, data_cpu_w, cmd_cpu_w);
20
21     logic[cache_line_size * BITS_IN_BYTE - 1:0] buff;
22     logic[cache_tag_size + cache_offset_size + cache_set_size - 1 : 0] test_addr = 19'b1111010101001001001;
23     logic[4 * BITS_IN_BYTE:0] test_payload = $random();
24     integer i;
25     longint timing;
26
27     task resetCache();
28         reset <= 1;
29         @(posedge clk);
30         reset <= 0;
31         @(negedge clk);
32     endtask
33
34     initial begin
35         begin : TEST_CACHE_HITS
36             resetCache();
37             driver.run_read(0, C1_READ8, buff, timing);
38             if (timing != cache_miss_time)
39                 $display("Cache read-miss timing error! Expected: %d, Real %d", cache_miss_time, timing);
40             for (i = 1; i < 32; ++i) begin
41                 driver.run_read(i, C1_READ8, buff, timing);
42                 if (i != 16 && timing != cache_hit_time) // i = 16 -- cache miss
43                     $display("Cache read-hit timing error! Expected: %d, Real %d", cache_hit_time, timing);
44             end
45             if (total_hits != 30)
46                 $display("Wrong cache hits! Expected: 30, Found: %d", total_hits);
47         end
48
49         begin : TEST_SINGLE_READ_WRITE_32
50             resetCache();
51             driver.run_write(test_addr, C1_WRITE32, test_payload, timing);
52             if (timing != cache_miss_time)
53                 $display("Cache write-miss timing error! Expected: %d, Real %d", cache_miss_time, timing);
54             driver.run_read(test_addr, C1_READ32, buff, timing);
55             if (buff[0 +: BITS_IN_BYTE * 4] != test_payload[0 +: BITS_IN_BYTE * 4]) begin
56                 $display("Cache correctness qword unit test failed, real: %b expected: %b",
57                     buff[0 +: BITS_IN_BYTE * 4], test_payload[0 +: BITS_IN_BYTE * 4]);
58             end
59             driver.run_write(test_addr, C1_WRITE32, test_payload, timing);
60             if (timing != cache_hit_time)
61                 $display("Cache write-hit timing error! Expected: %d, Real %d", cache_hit_time, timing);
62         end
63
64         begin : TEST_SINGLE_READ_WRITE_16
65             resetCache();
66             driver.run_write(test_addr, C1_WRITE16, test_payload, timing);
67             if (timing != cache_miss_time)

```

```

68     $display("Cache write-miss timing error! Expected: %d, Real %d", cache_miss_time, timing);
69     driver.run_read(test_addr, C1_READ16, buff, timing);
70     if (buff[0 +: BITS_IN_BYTE * 2] != test_payload[0 +: BITS_IN_BYTE * 2]) begin
71         $display("Cache correctness dword unit test failed, real: %b expected: %b",
72             buff[0 +: BITS_IN_BYTE * 2], test_payload[0 +: BITS_IN_BYTE * 2]);
73     end
74     driver.run_write(test_addr, C1_WRITE16, test_payload, timing);
75     if (timing != cache_hit_time)
76         $display("Cache write-hit timing error! Expected: %d, Real %d", cache_hit_time, timing);
77     end
78
79     begin : TEST_SINGLE_READ_WRITE_8
80         resetCache();
81         driver.run_write(test_addr, C1_WRITE8, test_payload, timing);
82         driver.run_read(test_addr, C1_READ8, buff, timing);
83         if (buff[0 +: BITS_IN_BYTE] != test_payload[0 +: BITS_IN_BYTE]) begin
84             $display("Cache correctness word unit test failed, real: %b expected: %b",
85                 buff[0 +: BITS_IN_BYTE], test_payload[0 +: BITS_IN_BYTE]);
86         end
87         driver.run_write(test_addr, C1_WRITE8, test_payload, timing);
88         if (timing != cache_hit_time)
89             $display("Cache write-hit timing error! Expected: %d, Real %d", cache_hit_time, timing);
90     end
91
92     $display("Finish cache testing");
93     $finish;
94 end
95
96 endmodule

```

Listing 8: cache_testbench.sv