

Mehrstufige Caches

Xuanqi Meng, Jeremias Rieser, Artem Bilovol

GRA SystemC Gruppe 192

2024

Inhalt I

Problemstellung

- Von-Neumann-Flaschenhals

Lösungsansätze und Optimierungen

- Cache-Architektur

- Assoziativität

- Direkt abgebildeter Cache

- Write-Through Cache

- Literaturrecherche

Korrektheit und Genauigkeit

- Korrektheit und Genauigkeit

- Korrektheit bei Schreiben und Lesen

- Anzahl von Hits und Misses

- Anzahl von Hits und Misses

- Auswirkungen

- Anzahl von Cycles

- Aligned vs unaligned accesses

Inhalt II

Schaltkreisanalyse

- Simulation

- L1 Control Unit

- L2 Control Unit

- Gatterberechnung

- Vergleich mit realen Systemen

Zusammenfassung und Ausblick

- Zusammenfassung

- Ausblick

Referenzen

Zusätzliches Material

Problemstellung

- ▶ Moderne Prozessoren sind deutlich schneller als der Hauptspeicher.
- ▶ Von Neumann-Flaschenhals: Die beschränkte Bandbreite und langsamen Datenübertragungsraten zwischen Prozessor und Hauptspeicher führen zu einem Leistungsengpass..
- ▶ Ziel: Untersuchung der Auswirkungen von mehrstufigen Caches auf Laufzeit und Latenz.

Von-Neumann-Flaschenhals

- ▶ Prozessoren können Daten schneller verarbeiten als diese aus dem Hauptspeicher gelesen werden können.
- ▶ Der Hauptspeicher stellt einen Flaschenhals dar, der die Gesamtleistung des Systems limitiert.
- ▶ Lösung: Einsatz von Caches als Puffer zwischen Prozessor und Hauptspeicher.

Lösungsansätze und Optimierungen

- ▶ Implementierung einer Cache-Simulation in SystemC und C++.
- ▶ Simulation eines direkt-assoziativen Caches.
- ▶ Untersuchung von Parametern wie Cachegröße, Assoziativität, Cachezeilenanzahl und Zeilengröße und deren Auswirkung auf die Zugriffszeiten.

Cache-Architektur

- ▶ Verwendung von L1- und L2-Caches zur Reduktion der Zugriffszeit [4]
- ▶ L1-Cache: kleine Größe aber geringe Latenz (wenige Zyklen)
- ▶ L2-Cache: größer als L1-Cache (hier: inklusive) aber höhere Latenz
- ▶ Background-memory: Hauptspeicher der Systems, höchste Latenz
- ▶ Inklusivität: alle Adressen, die in L1 speichern auch in L2 aber nicht umgekehrt.

Assoziativität

- ▶ Direkt abgebildeter Cache: Jeder Speicherblock verwaltet eine Speicheradresse (auch einfach assoziativ genannt)
- ▶ Alternativ: höhere Assoziativität erlaubt effizientere Cache-Nutzung aber höhere Kosten

Direkt abgebildeter Cache

- ▶ Jeder Speicherblock hat genau eine mögliche Position im Cache.
- ▶ Einfacher Aufbau und schnelle Zugriffszeiten.
- ▶ Nachteile: Höhere Konfliktwahrscheinlichkeit und niedrigere Trefferquote im Vergleich zu assoziativen Caches.

Write-Through Cache

- ▶ Write-Through: Daten werden gleichzeitig in den Cache und den Hauptspeicher geschrieben.
- ▶ Vorteile: Daten im Hauptspeicher sind immer aktuell.
- ▶ Nachteile: kumulative Latenz bei Schreiboperationen.

Daten- und Adressbus

- ▶ Datenbus und Adressbus sind 4 Byte breit.
- ▶ Der Hauptspeicher ist jedoch byte-adressiert.

- ▶ Übliche Cachegrößen: L1 (32KB), L2 (256KB) [5]
- ▶ Latenzen: L1 (3 - 5 Zyklen), L2 (10 - 20 Zyklen) [1],[5]
- ▶ Vergleich der eigenen Implementierung mit den Ergebnissen

Typedefs

```
1 typedef struct Request {
2     uint32_t addr;
3     uint32_t data;
4     int we;    // 0 read 1 write
5 } Request;
6
7 typedef struct Result {
8     size_t cycles;
9     size_t misses;
10    size_t hits;
11    size_t primitiveGateCount;
12 } Result;
13
14 typedef struct CacheLine {
15     int tag;
16     uint8_t *bytes;
17     int empty;
18 } CacheLine;
```

Listing 1: types.h

CACHEL1.h

```
1  SC_MODULE(CACHEL1){
2      ...
3      CacheLine *internal;
4
5      sc_in<bool> clk;
6      sc_in<bool> requestIncoming;
7      sc_in<sc_bv<32>> inputData;
8      sc_in<sc_bv<32>> address;
9      sc_in<bool> rw;
10     sc_out<bool> ready;
11     sc_out<sc_bv<32>> outputData;
12     sc_out<bool> isWriteThrough;
13     ...
14     void run();
15     void write();
16     void read();
17     void writeThrough(int index, int address);
18     int loadFromL2(int address);
19     ...
20 };
```

Listing 2: CACHEL1.h - Cache Level 1 Header

CACHEL1.h

```
1 SC_CTOR(CACHEL1);
2 CACHEL1(sc_module_name name, int latency, int cacheLines, int cacheLineSize)
3 : sc_module(name) {
4     /*...*/
5     offsetLength = (int)(log(cacheLineSize) / log(2));
6     indexLength = (int)(log(cacheLines) / log(2));
7     tagbits = 32 - offsetLength - indexLength;
8     tagOffset = 32 - tagbits;
9     indexOffset = 32 - tagbits - indexLength;
10
11     internal = new CacheLine[cacheLines];
12     for (int i = 0; i < cacheLines; i++) {
13         internal[i].bytes = (uint8_t*)calloc(cacheLineSize, sizeof(uint8_t));
14         internal[i].empty = 1;
15     }
16
17     SC_THREAD(run);
18     sensitive << clk.pos() << requestIncoming;
19 }
```

Memory.h

```
1 SC_MODULE(MEMORY){
2     int latency;
3     sc_in<bool> clk;
4
5     sc_in<bool> requestIncoming;
6     sc_in<bool> rw;
7     sc_in<sc_bv<32>> addr;
8     sc_in<sc_bv<32>> rData; // the incoming data from L2
9
10    sc_out<bool> ready;
11    sc_out<sc_bv<32>> wData; // output data to L2
12
13    std::unordered_map<uint32_t,uint8_t> internal;
14
15    SC_CTOR(MEMORY);
16    MEMORY(sc_module_name name, int latency) : sc_module(name) {
17        this->latency = latency;
18        SC_THREAD(run);
19        sensitive << clk.pos() << requestIncoming;
20    }
21    void run();
```


Überblick über die Simulationsfunktion

```
1 Result run_simulation(int cycles, unsigned l1CacheLines,
2                       unsigned l2CacheLines,
3                       unsigned cacheLineSize, unsigned l1CacheLatency,
4                       unsigned l2CacheLatency, unsigned memoryLatency,
5                       size_t numRequests, struct Request requests[],
6                       const char *tracefile) {
7
8     MEMORY memory("mem",memoryLatency);
9     CACHEL2 l2Cache("l2",l2CacheLatency,l2CacheLines,cacheLineSize);
10    CACHEL1 l1Cache("l1",l1CacheLatency,l1CacheLines,cacheLineSize);
11
12    l1Cache.l2 = l2Cache.internal;
13    l2Cache.l1 = l1Cache.internal;
14
15    };
```

Weiter in der Simulationsfunktion

```
1 // one tick for initialization
2   clk.write(true);
3   sc_start(1,SC_NS);
4   clk.write(false);
5   sc_start(1,SC_NS);
6
7   bool lastR = false;
8   int i;
9   for(i = 0; i < cycles ; i++){
10     /*...*/ // pr fe ob alle Module bereit sind und schicke die Anfrage
11     clk.write(true);
12     sc_start(1,SC_NS);
13     requestToL1.write(false);
14     clk.write(false);
15     sc_start(1,SC_NS);
16   }
```

Korrektheit und Genauigkeit

- ▶ Korrektheit bei Schreiben und Lesen: Es soll kein Unterschied zwischen Zugriffen mit und ohne Cache bestehen.
 - ▶ Beispiel für erfolgreiche Speicherung von Daten im Speicher und deren korrekte Darstellung beim Lesen.
 - ▶ Automatischer Vergleich von Daten aus Cache und ohne Cache bei 1000 zufällig generierten Anfragen (Kommandozeilenargumente werden sowohl manuell als auch automatisch erzeugt).
- ▶ Anzahl von Hits und Misses:
 - ▶ Codeanalyse
 - ▶ Beispiel
- ▶ Anzahl von Cycles:
 - ▶ Bestimmung der Anzahl von Zyklen beim Lesen und Schreiben
 - ▶ Analyse von Beispielen

Korrektheit bei Schreiben und Lesen

Ein Beispiel: CachelineSize = 64

Parameter:

--l1-lines 2 --l2-lines 2 e.csv

Eingabe:

```
W,0x0000003e,0x12345678
W,0x000000be,0x87654321
R,0x0000003d,
R,0x000000bd,
R,0x0000003c,
```

```
Request 0:
    addr: 0x0000003E
    data: 0x12345678
    we: 1
Request 1:
    addr: 0x000000BE
    data: 0x87654321
    we: 1
Request 2:
    addr: 0x0000003D
    data: 0x00123456
    we: 0
Request 3:
    addr: 0x000000BD
    data: 0x00876543
    we: 0
# ...
```

Korrektheit bei Schreiben und Lesen (Fort.)

Ein Beispiel: CachelineSize = 64

Parameter: --l1-lines 2 --l2-lines 2
e.csv

Eingabe:

```
W,0x0000003e,0x12345678
W,0x000000be,0x87654321
R,0x0000003d,
R,0x000000bd,
R,0x0000003c,
```

Ausgabe:

```
# ...
Request 3:
    addr: 0x000000BD
    data: 0x00876543
    we: 0
Request 4:
    addr: 0x0000003C
    data: 0x00001234
    we: 0
Result:
    cycles: 58570
    misses: 14
    hits: 0
    primitiveGateCount: 42649
```

Korrektheit bei Schreiben und Lesen (Fort.)

- ▶ Erste Stufe nur mit dem Memory-modul.
- ▶ Kein Overhead bzw. keine Fehleranfälligkeit durch mehrstufiges Cachen
- ▶ Als input zufällig erzeugte .csv (1000 Zeilen) und vergleich der Ergebnisse
- ▶ Test ist natürlich nur sinnvoll, wenn sich die Adressen wiederholen bzw. nahe aneinander liegen (überschneidung)
- ▶ gesichert durch simulation von 1000 Adressen aus einem Adresspool von Größe 500
- ▶ eine Testeingabe mit parameter -l1-lines 2 -l2-lines 2:

```
W,0x00000581,0x66666666
W,0x00000440,0x33333333
R,0x00000070,
R,0x00001CE1,
#...
W,0x00000901,0x22222222
W,0x00001B13,0x55555555
R,0x00001C60,
```

Korrektheit bei Schreiben und Lesen(Fortsetzung)

- ▶ Die Ergebnisse werden dann durch eine .sh verglichen

```
1  for INPUTFILE in *.csv; do
2      if [ -f "$INPUTFILE" ]; then
3          ./project -args ... > to_test.txt
4          ./test/project -args ... > direct_memory.txt # this is only memory
5          if diff -q to_test.txt direct_memory.txt > /dev/null; then
6              echo "No difference between the output with and without cache."
7          else
8              echo "There is a difference, something is wrong."
9          fi
10     fi
11 done
```

- ▶ Ausgabe von .sh zeigt die Korrektheit

```
Checking cache_accesses.csv ...
      No difference between the output with and without cache.
The test is over.
```

Korrektheit bei Schreiben und Lesen(Fortsetzung)

- ▶ zusätzliches Testen durch .sh mit zufällig generierten Kommandozeilenargumenten

```
1  ...
2
3  CYCLES=$(( (RANDOM * RANDOM) % 100000000 ))
4
5  LINESIZE=$(( 2 ** (2 + RANDOM % 13) ))
6
7  L2LINES=$(( 2 ** (2 + RANDOM % 13) ))
8
9  L1LINES=$(( 2 ** (2 + RANDOM % 13) ))
10
11 ...
```

- ▶ Ausgabe von .sh zeigt die Korrektheit

```
The test is over.
No problems found.
```


Anzahl von Hits und Misses

► Codeanalyse (Mechanismus von Hit und Miss)

```
1  index = ifExist(t_tmp,i_tmp);
2  if(index==-1){ // in the cache there are no such information
3      miss++;
4      hit = false;
5      index = loadFromL2(address.read().to_int());
6  }
7  if(hit){
8      hits++;
9  }
```

Listing 3: L1Cache.cpp

- ifExist(tag,index) sucht die zugreifende Zeile
- Ein Ähnlicher Mechanismus kommt bei L2 zum einsatz
- Bei einem zeileübergreifenden Zugriff gilt es nur als Hit, wenn beide Zeilen gefunden werden können

Berechnung von Hits und Misses (Fortsetzung)

► Berechnung von Hit und Miss

- Die Hits und Misses von beiden L1 und L2 werden zusammen gezählt
- Beim Read (Zugriff in 1 Zeile):
 - L1 Hit = 1 Hit
 - L1 Miss + L2 Hit = 1 Miss + 1 Hit
 - L1 Miss + L2 Miss = 2 Miss
- (Zugriff in 2 Zeilen)
 - L1 Hit = 1 Hit
 - L1 Miss(1 Zeile) + L2 Hit = 1 Miss + 1 Hit
 - L1 Miss(1 Zeile) + L2 Miss = 2 Miss
 - L1 Miss(2 Zeile) + L2 Hit + L2 Hit = 1 Miss + 2 Hits
 - L1 Miss(2 Zeile) + L2 Miss + L2 Hit = 2 Miss + 1 Hit (2 mal Anfrage zur L2)
 - L1 Miss(2 Zeile) + L2 Miss + L2 Miss = 3 Miss
- Beim Write ist der einzige Unterschied: L1 Hit + L2 Hit = 2 Hit
- Schreiben findet in allen Modulen statt

Berechnung von Hits und Misses (Beispiel)

► Berechnung von Hit und Miss

Ein Beispiel: CachelineSize = 4

Parameter: --l1-lines 2 --l2-lines 4

```
1 W,0x0000003e,0x12345678 #Unaligned L1 miss + L2 miss + L2 miss
2 #2 Zeile in L1 und L2 geladt
3 W,0x0000003f,0x87654321 #Unaligned L1 hit + L2 hit
4 #2 Zeile existert in L1 und L2
5 W,0x0000003c,0x56789123 #aligned L1 hit + L2 hit
6 #1 Zeile existiert in L1 und L2
7 W,0x00000037,0x12345678 #unaligned L1 miss + L2 miss + L2 miss
8 #2 Zeile in L1 und L2 geladt
9 W,0x0000003d,0x12345678 #unaligned L1 miss + L2 hit + L2 hit
10 #2 Zeile in L1 geladt, ist aber in L2 vorhanden
11 R,0x000000f4, #aligned L1 miss + L2 miss
12 #1 Zeile wird in beide L1 und L2 geladt
13 R,0x0000003d, #unaligned L1 miss(1 Zeile) + L2 hit
14 #1 Zeile existiert nicht in L2, wird von Memory geladt
```

Berechnung von Hits und Misses (Hit Log)

► Compile mit Flag HIT_LOG:

```
l1 received <write> with addr:0000003e and data:12345678
l1 miss by writing first line with addr: 0000003e detected, sending
  signal to next level
l2 miss by reading with addr: 0000003e detected, sending signal to
  next level
l1 miss by writing second line with addr: 00000041 detected, sending
  signal to next level
l2 miss by reading with addr: 00000041 detected, sending signal to
  next level
l1 miss by writing by 2 lines with addr:0000003e
# ...
l1 received <read> with addr:000000f4
l1 miss by reading with addr: 000000f4 detected, sending signal to
  next level
l2 miss by reading with addr: 000000f4 detected, sending signal to
  next level
l1 received <read> with addr:0000003d
l1 miss by reading first line with addr: 0000003d detected, sending
  signal to next level
```

Berechnung von Hits und Misses (Ergebnis)

► Ausgegebene Anzahl an Hits und Misses

Result:

cycles: 177

misses: 10

hits: 7

primitiveGateCount: 4369

Auswirkungen von Parametern

- Zu kleine Cachelines nutzen die Lokalität der Inputs nicht richtig aus, also mehr misses, gerade bei unaligned-access [6]

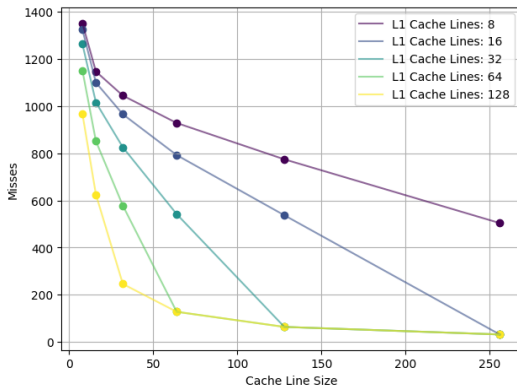
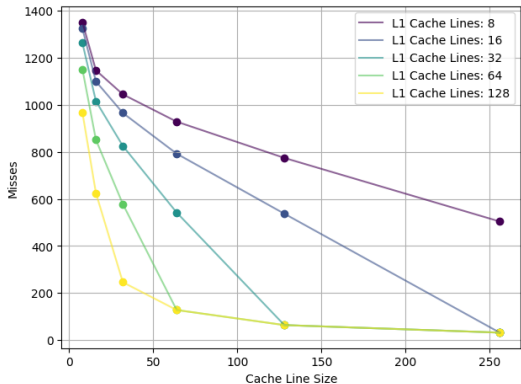


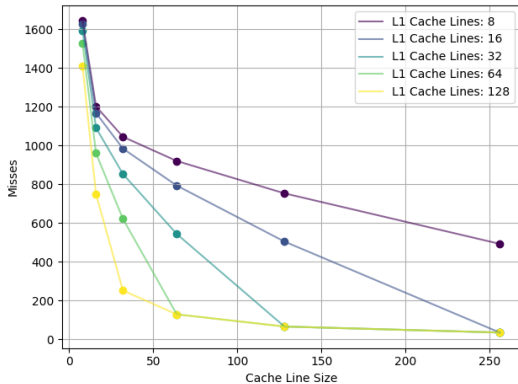
Figure: Test with 1000 requests (aligned)

Auswirkungen von Parametern (Forts.)

- ▶ Aligned accesses haben durchschnittlich mehr misses als aligned accesses, vorallem bei kleinen Cachelines.



Test with 1000 requests (aligned)



Test with 1000 requests (unaligned)

Anzahl von Cycles

- ▶ Anzahl von Cycles sind deterministisch je nach Verhalten des Inputs
- ▶ Für einen Datenzugriff nehmen wir die Latenzen in den verschiedenen Modulen
- ▶ In modernen Caches ist es deutlich schneller, nur den Tag anstatt von gesamten Daten zu verändern.

Anzahl von Cycles

- ▶ Anzahl von Cycles bei verschiedenen Verhalten

- ▶ Bei Read

L1 Hit = L1 Latency

L1 Miss + L2 Hit = L1 Latency + L2 Latency

L1 Miss + L2 Miss = L1 Latency + L2 Latency + cachelineSize/4* Memory Latency

(Zugriff in 2 Zeilen)

L1 Hit + L1 Hit = 2*L1 Latency

L1 Miss + L1 Hit + L2 Hit = 2* L1 Latency + L2 Latency

L1 Miss + L1 Hit + L2 Miss = 2*L1 Latency + L2 Latency + cachelineSize/4* Memory Latency

L1 Miss + L1 Miss + L2 Hit + L2 Hit = 2*L1 Latency + 2* L2 Latency

L1 Miss + L1 Miss + L2 Miss + L2 Hit = 2*L1 Latency + 2* L2 Latency + cachelineSize/4* Memory Latency

L1 Miss + L1 Miss + L2 Miss + L2 Miss = 2*L1 Latency + 2* L2 Latency + 2*(cachelineSize/4* Memory Latency)

Anzahl von Cycles(Fortsetzung)

- ▶ Anzahl von Cycles bei verschiedenen Verhalten
- ▶ Bei Write(Through) — Write-Allocate wird angepasst
 - $L1 \text{ Hit} = \text{Memory Latency}$
 - $L1 \text{ Miss} + L2 \text{ Hit} = L2 \text{ Latency} + \text{Memory Latency}$
 - $L1 \text{ Miss} + L2 \text{ Miss} = L2 \text{ Latency} + \text{cachelineSize}/4 * \text{Memory Latency} + \text{Memory Latency}$
 - ...
- ▶ Alle L1 Latenzen werden in einer Memory Latenz mitgezählt, solange $\text{Memory Latency} > 2 * L2 \text{ Latency} > 2 * L1 \text{ Latency}$
Ansonsten: Memory Latency wird durch die größte der oberen drei ersetzt.

Anzahl von Cycles (Beispiel)

► Berechnung von Hit und Miss

Ein Beispiel: CachelineSize = 8

Parameter: --l1-lines 2 --l2-lines 4 --l1-latency 3 --l2-latency 5
--memory-latency 12

```
1 W,0x0000003e,0x12345678 #Unaligned L1 miss + L1 miss + L2 miss + L2
   miss
2                               # 2*5 + 2*(8/4)*12 + 12 = 70
3 W,0x0000003f,0x87654321 #Unaligned L1 hit + L1 hit + L2 hit + L2 hit
4                               # 12
5 W,0x00000038,0x11111111 #aligned L1 hit + L2 hit
6                               # 12
7 W,0x0000004f,0x12345566 #unaligned L1 miss + L1 miss + L2 miss + L2
   miss
8                               # 2*5 + 2*(8/4)*12 + 12 = 70
9 W,0x0000003d,0xdddddddd #unaligned L1 miss + L1 miss + L2 hit + L2
   hit
10                               # 2*5 + 12 = 22
11 #...
```

Anzahl von Cycles (Beispiel—Fortsetzung)

► Berechnung von Hit und Miss

Ein Beispiel: CachelineSize = 8

Parameter: --l1-lines 2 --l2-lines 4 --l1-latency 3 --l2-latency 5
--memory-latency 12

```
1 #...
2 R,0x000000ec,          #aligned    L1 miss + L2 miss
3                   # 3 + 5 + (8/4)*12 = 32
4 R,0x0000003d,          #unaligned L1 miss + L1 hit + L2 hit
5                   # 2*3 + 5 = 11
6                   # insgesamt : 229
```

Anzahl von Cycles (Ergebnis)

► Ergebnis

Result:

```
cycles: 237  
misses: 10  
hits: 7  
primitiveGateCount: 6991
```

► Overhead

1. Bei jedem Ready von L1 wird noch ein Cycle gewartet bis zur nächsten Anfrage
also $7 \text{ Anfragen} + 229 = 236$
2. Nach dem Ready von allen Module wird noch ein Cycle gewartet bis zum Ende
des Programms
also $236 + 1 = 237$

► die Anzahl an Cycles sind also deterministisch gegenüber der Operation

Aligned vs Unaligned access

- ▶ Zugriffe sollen auch mit Offset möglich sein
- ▶ Nachteil bei der Performanz [3]:
- ▶ Generell brauchen unaligned-accesses in manchen Fällen (abhängig von Größe) Zugriff auf 2 Cachelines - Damit verdoppelt sich die Zugriffslatenz bei Misses
- ▶ Dies ist im Test korrekt wiedergegeben

Schaltkreis

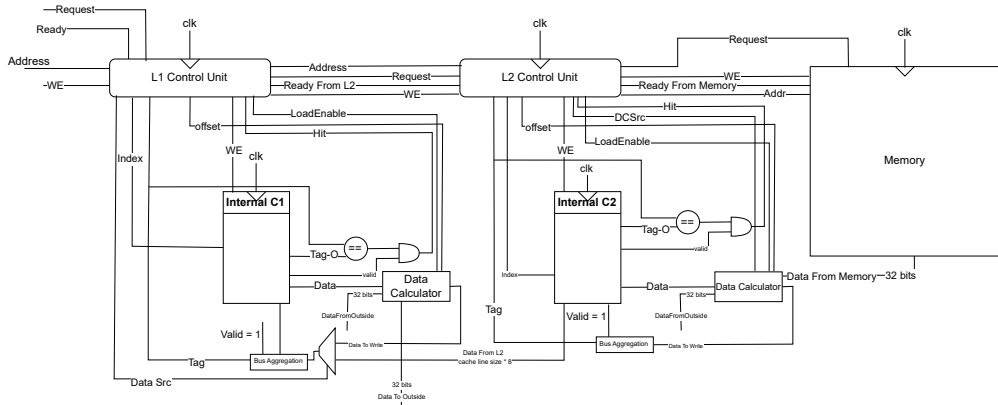


Figure: Schaltkreis - Simulation

L1 Control Unit

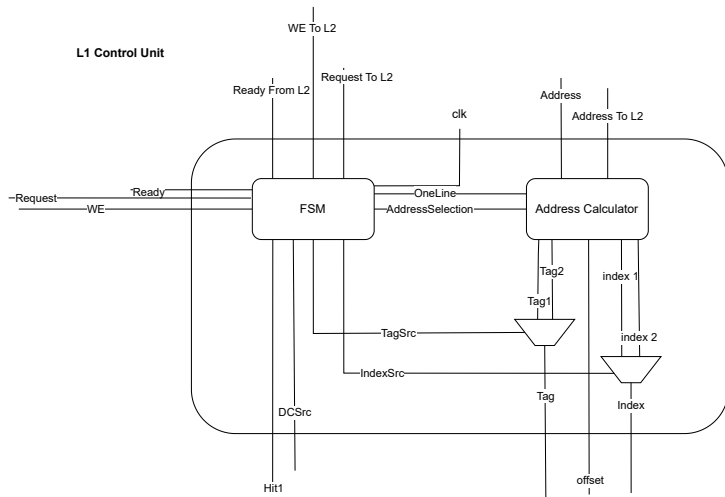


Figure: L1 Control Unit

L1 Control Unit - FSM

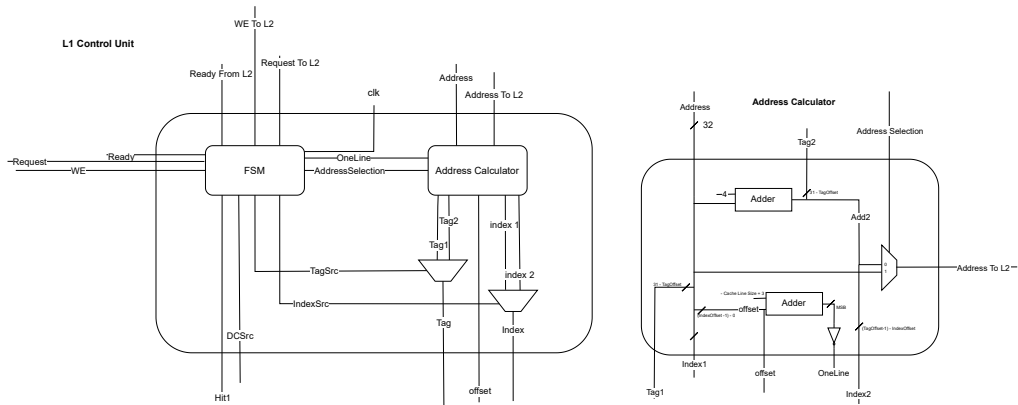


Figure: L1 Control unit and Address calculator

L1 Control Unit - FSM

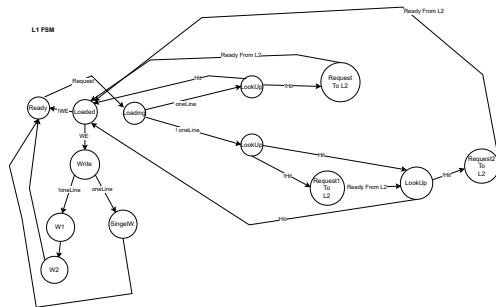
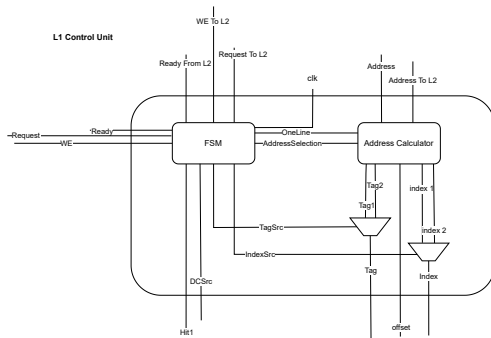


Figure: L1 Control unit and FSM

L2 Control Unit

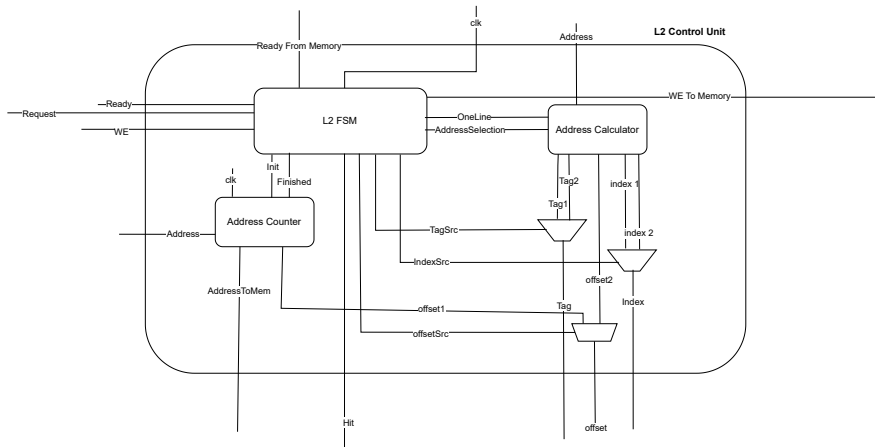


Figure: L2 Control Unit

L2 Control Unit - FSM

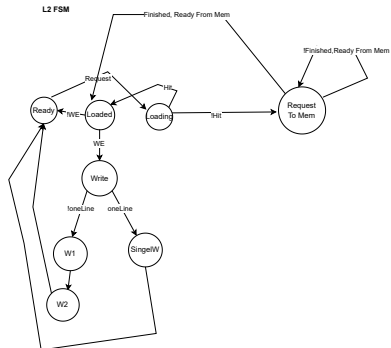
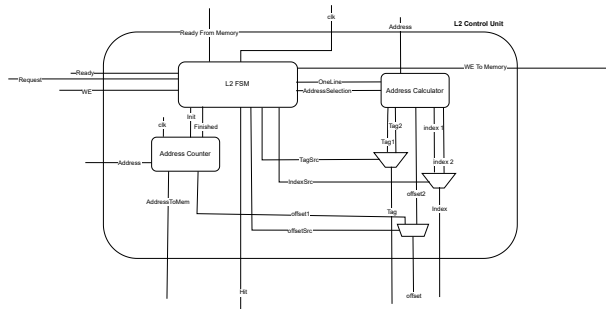


Figure: L2 Control unit and FSM

L2 Control Unit - Address Counter

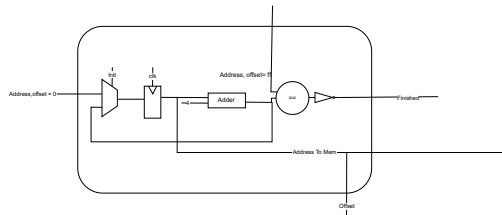
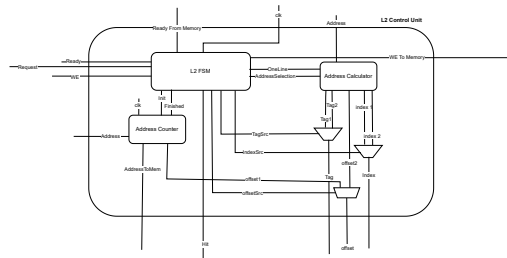


Figure: L2 Control unit and Address Counter

Gatterberechnung L1

► L1:

Speicherung: $(\text{Tagbits} + \text{valid} + \text{CachelineSize} * 8) * \text{Cachlines}$

Tag Comparator: $\text{TagLength} + \text{TagLength} - 1 \text{ (XOR und OR)} + 1 \text{ (AND)}$

CachelineSize*8-bit-2-to-1 MUX : $\text{CachelineSize} * 8 * 4$

Control Unit:

2-to-1 MUX für Tag: $\text{TagLength} * 4$

2-to-1 Mux für Index: $\text{IndexLength} * 4$

FSM:

Zustandspeicherung: $\log_2(13) = 4 \text{ (bits)}$ D-Flip-Flop, Gatter = $4 * 5 = 20$

Zustandübergang: $(5 \text{ Eingabe} + 4 \text{ bits}) * 4 \text{ bits} = 36$

Ausgabe: $4 \text{ bits} * 7 \text{ Ausgabe} = 28$

Gatterberechnung L1 (Forts.)



L1:

Address Calculator:

32-bit Adder : $(2 \text{ And} + 1 \text{ Or} + 2 \text{ XOR}) * 32 = 160$

offset Adder : $(2 \text{ And} + 1 \text{ Or} + 2 \text{ XOR}) * \text{offsetLength}$

32-bit-2-to-1 MUX: $32 * 4 = 128$

Data Calculator:

Speicherung: $(2 * \text{CachelineSize} * 8)$ D-Latch(jeweils 5 Gatter)

1-to-2 Decoder: 1 Not

$2 * 1\text{-bit-2-to-1 MUX}: 2 * 4$

$\text{CachelineSize} * 8\text{-bit-2-to-1 MUX}: \text{CachelineSize} * 8 * 4$

$\text{CachelineSize} * 32\text{-bit-2-to-1 MUX} : \text{CachelineSize} * 32 * 4$

$\text{CachelineSize} * 32\text{-bit-2-to-1 DEMUX} : \text{CachelineSize} * 32 * 3$

Gatterberechnung L2

► L2:

Speicherung: $(\text{Tagbits} + \text{valid} + \text{CachelineSize} * 8) * \text{Cachelines}$

Tag Comparator: $\text{TagLength} + \text{TagLength} - 1$ (XOR und OR) + 1 (AND)

Control Unit:

2-to-1 MUX für Tag: $\text{TagLength} * 4$

2-to-1 MUX für Index: $\text{IndexLength} * 4$

2-to-1 MUX für Offset : $\text{OffsetLength} * 4$

FSM:

Zustandspeicherung: $\log_2(8) = 3$ (bits) D-Flip-Flop, Gatter = $3 * 5 = 15$

Zustandübergang: $(6 \text{ Eingabe} + 3 \text{ bits}) * 3 \text{ bits} = 27$

Ausgabe: $3 \text{ bits} * 7 \text{ Ausgabe} = 21$

Gatterberechnung L2 (Forts.)

► L2:

Address Counter:

32-bit-2-to-1 MUX : $32 * 4 = 128$

32-bit D-Flip-Flop = $32 * 5 = 160$

32-bit Adder = 160

32-bit Comparator = $32 \text{ XOR} + 31 \text{ OR} = 63$

1 Inverter = 1

Address Calculator:

gleich wie bei L1

Data Calculator:

Speicherung: $(2 * \text{CachelineSize} * 8)$ D-Latch(jeweils 5 Gatter)

1-to-2 Decoder: 1 Not

$2 * 1\text{-bit-2-to-1 MUX}$: $2 * 4$

$\text{CachelineSize} * 8\text{-bit-2-to-1 MUX}$: $\text{CachelineSize} * 8 * 4$

$\text{CachelineSize} * 32\text{-bit-2-to-1 MUX}$: $\text{CachelineSize} * 32 * 4$

Kein Demux nötig, da wir keine 32-bit Daten auslesen sollen

Vergleich mit realen Systemen

- ▶ Bei vielen SRAM-Chips ist neben der Speicherung häufig auch Funktionalität wie FSM, Adressverarbeitung und Datenverarbeitung eingebaut, die mit unseren Komponenten vergleichbar ist.
- ▶ z.B. CY7C1380D von CYPRESS [2]

Zusammenfassung

- ▶ Das Ergebnis dieses Projekts ist eine Implementierung eines zweistufigen Caches, der als Werkzeug zur Messung des Nutzens von Cache-Speicher dient.
- ▶ Im Verlauf des Projekts wurden mehrere Standardprobleme analysiert und gelöst, insbesondere:
 - ▶ Abstraktion von Hardwarekomponenten in der Simulation
 - ▶ Synchronisation zwischen Modulen
 - ▶ Datenausrichtung und Bytereihenfolge
 - ▶ Softwaretests
- ▶ Ein weiterer wichtiger Teil des Projekts war die Strukturierung und das Kompilieren, insbesondere unter Berücksichtigung der Umgebung, die nicht vom Team kontrolliert wird.

Basierend auf den Ergebnissen des Projekts werden folgende Bereiche für mögliche Verbesserungen in Betracht gezogen:

- ▶ Hinzufügen von L3-Cache und externem Cache zur Simulation
- ▶ Multiprocessing und unterschiedliche Implementationen von L1-Cache (insbesondere Cache-Fairness)
- ▶ I-Cache und D-Cache
- ▶ Testen aligned vs unaligned Offset-Zugriffe
- ▶ Simulation von anderen Arten von Cache (exklusive, andere Assozitivität usw.)

- [1] colin scott colin. *Numbers Every Programmer Should Know By Year*. 2021. URL: https://colin-scott.github.io/personal_website/research/interactive_latency.html (visited on 07/21/2024).
- [2] Cypress Semiconductor Corporation. *18-Mbit (512K × 36/1M × 18) Pipelined SRAM*. 2016. URL: https://www.mouser.com/datasheet/2/100/CYPR_S_A0002300802_1-2540696.pdf.
- [3] Daniel Drake. *Unaligned Memory Accesses — The Linux Kernel documentation*. URL: <https://www.kernel.org/doc/html/latest/core-api/unaligned-memory-access.html> (visited on 07/21/2024).
- [4] Christoph Ersfeld. “Organisation von Caches”. de. In: (2002).
- [5] “Intel® 64 and IA-32 Architectures Optimization Reference Manual: Volume 1”. en. In: (2024).
- [6] Viraj Kumar. *CS232: Computer Architecture II*. 2010. URL: <https://courses.grainger.illinois.edu/cs232/sp2010/lectures/L17.pdf> (visited on 07/21/2024).

Daten Berechner (Einlesen und Auslesen mit Cache)

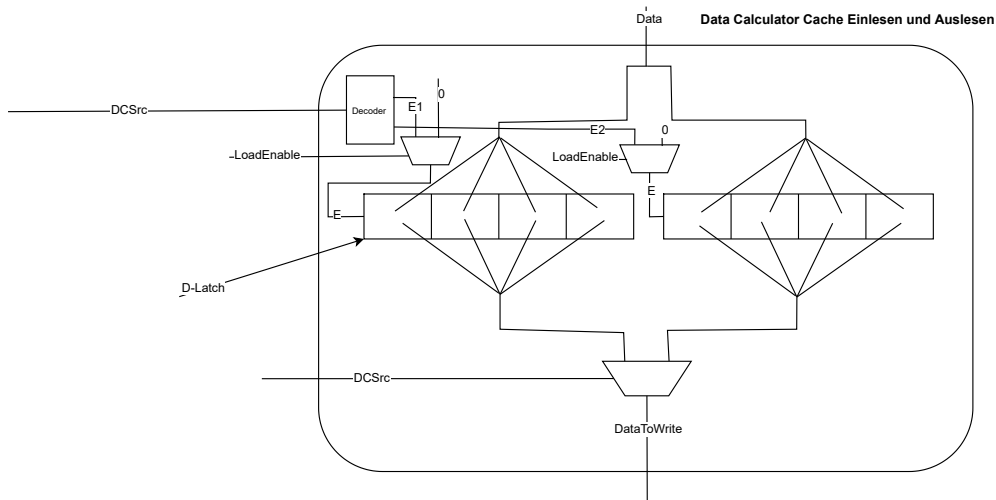


Figure: Einlesen und Auslesen mit Cache

Daten Berechner (Einlesen mit 32 bits Daten)

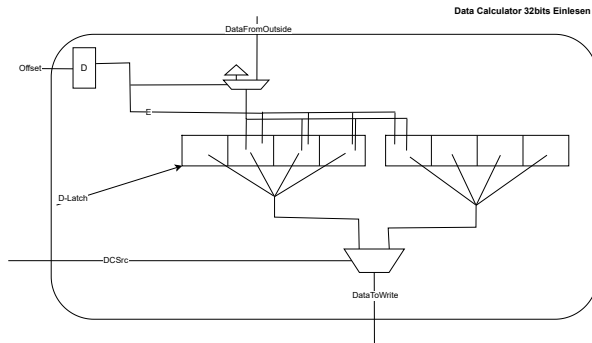


Figure: Einlesen mit 32 bits Daten

Daten Berechner (Auslesen mit 32 bits Daten)

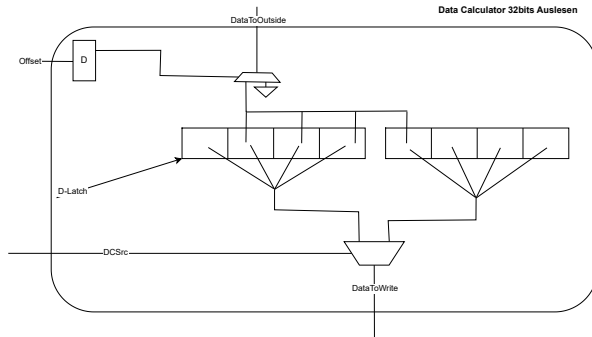


Figure: Auslesen mit 32 bits Daten