Introduction
○○

Addition
○○○○○○○○○○

Multiplication
○○○○○○○○○○○○

Division
○○○○○○○○○○

Conclusion
○○

# Efficient Big Integer Arithmetic Using GPGPU

With focus on implementations of exact addition, division, and multiplication in CUDA C++ and Futhark

Thorbjørn Bülow Bringgaard

University of Copenhagen
Department of Computer Science

June 18, 2024

# Outline

1. Introduction

2. Addition

3. Multiplication

4. Division

5. Conclusion

# Introduction
## Big Integers

### Positional number system

An integer $u \in \mathbb{N}$ can be expressed in base $B \in \mathbb{N} > 1$ with $m$ digits $u_{i \in \{0,1,\dots,m-1\}} \in \{0, 1, \dots, B-1\}$ by the sum:

$$u = \sum_{i=0}^{m-1} u_i \cdot B^i \tag{1}$$

E.g. the number 256 in base $B = 10$ is $[6, 5, 2]$.

Big integers in the positional number system maps to an array of unsigned machine words.

E.g. the number 4294967298 in base $B = 2^{32}$ is $[2, 1]$.

# Introduction
## Big Integers

### Positional number system

An integer $u \in \mathbb{N}$ can be expressed in base $B \in \mathbb{N} > 1$ with $m$ digits $u_{i \in \{0,1,\ldots,m-1\}} \in \{0, 1, \ldots, B-1\}$ by the sum:

$$u = \sum_{i=0}^{m-1} u_i \cdot B^i \tag{1}$$

E.g. the number 256 in base $B = 10$ is $[6, 5, 2]$.

Big integers in the positional number system maps to an array of unsigned machine words.

E.g. the number 4294967298 in base $B = 2^{32}$ is $[2, 1]$.

## Introduction
### Big Integers

#### Positional number system

An integer $u \in \mathbb{N}$ can be expressed in base $B \in \mathbb{N} > 1$ with $m$ digits
$u_{i \in \{0,1,\ldots,m-1\}} \in \{0, 1, \ldots, B-1\}$ by the sum:

$$u = \sum_{i=0}^{m-1} u_i \cdot B^i \tag{1}$$

E.g. the number 256 in base $B = 10$ is $[6, 5, 2]$.

Big integers in the positional number system maps to an array of unsigned machine words.

E.g. the number 4294967298 in base $B = 2^{32}$ is $[2, 1]$.

## Introduction
### Big Integers

#### Positional number system

An integer $u \in \mathbb{N}$ can be expressed in base $B \in \mathbb{N} > 1$ with $m$ digits
$u_{i \in \{0,1,\ldots,m-1\}} \in \{0, 1, \ldots, B-1\}$ by the sum:

$$u = \sum_{i=0}^{m-1} u_i \cdot B^i \tag{1}$$

E.g. the number 256 in base $B = 10$ is $[6, 5, 2]$.

Big integers in the positional number system maps to an array of
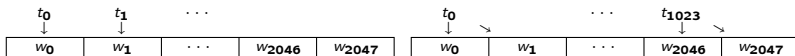unsigned machine words.

E.g. the number 4294967298 in base $B = 2^{32}$ is $[2, 1]$.

## Introduction
Implementation

The idea is to process integers at CUDA block-level, requiring:

- Sequentialization factor.

| $t_0$ | $t_1$ | $\cdots$ | | |
|-------|-------|----------|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

| $t_0$ | | $\cdots$ | $t_{1023}$ | |
|-------|---|----------|------------|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Coalesced transactions to global memory.

| $t_0$ | | $\cdots$ | $t_{1023}$ | |
|-------|---|----------|------------|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

| $t_0$ | $t_1$ | $\cdots$ | $t_{1022}$ | $t_{1023}$ |
|-------|-------|----------|------------|------------|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Multiple instances per block.

| $t_0$ | $t_1$ | $\cdots$ | $t_{14}$ | $t_{15}$ |
|-------|-------|----------|----------|----------|
| $w_0$ | $w_1$ | $\cdots$ | $w_{14}$ | $w_{15}$ |

| $t_0$ | $t_1$ | $\cdots$ | $t_{30}$ | $t_{31}$ |
|-------|-------|----------|----------|----------|
| $w_0$ | $w_1$ | $\cdots$ | $u_{14}$ | $u_{15}$ |

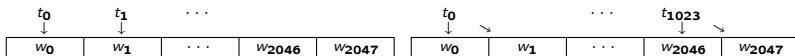# Introduction
Implementation

The idea is to process integers at CUDA block-level, requiring:

- Sequentialization factor.

| $t_0$ | $t_1$ | $\cdots$ | | | $t_0$ | | $\cdots$ | $t_{1023}$ | |
|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ | $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Coalesced transactions to global memory.

| $t_0$ | | $\cdots$ | $t_{1023}$ | | $t_0$ | $t_1$ | $\cdots$ | $t_{1022}$ | $t_{1023}$ |
|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ | $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Multiple instances per block.

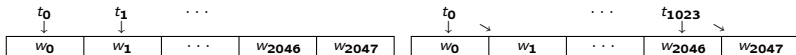| $t_0$ | $t_1$ | $\cdots$ | $t_{14}$ | $t_{15}$ | $t_0$ | $t_1$ | $\cdots$ | $t_{30}$ | $t_{31}$ |
|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{14}$ | $w_{15}$ | $w_0$ | $w_1$ | $\cdots$ | $u_{14}$ | $u_{15}$ |

## Introduction
Implementation

The idea is to process integers at CUDA block-level, requiring:

- Sequentialization factor.

| $t_0$ | $t_1$ | $\cdots$ | | $t_{1023}$ | |
|---|---|---|---|---|---|
| $\downarrow$ | $\downarrow$ | | | $\downarrow$ $\searrow$ | |
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

| $t_0$ | $\searrow$ | $\cdots$ | | $t_{1023}$ $\searrow$ | |
|---|---|---|---|---|---|
| $\downarrow$ | | | | $\downarrow$ | |
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Coalesced transactions to global memory.

| $t_0$ | $\searrow$ | $\cdots$ | $t_{1023}$ $\searrow$ | |
|---|---|---|---|---|
| $\downarrow$ | | | $\downarrow$ | |
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

| $t_0$ | $t_1$ | $\cdots$ | $t_{1022}$ | $t_{1023}$ |
|---|---|---|---|---|
| $\downarrow$ | $\downarrow$ | | $\downarrow$ | $\downarrow$ |
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Multiple instances per block.

| $t_0$ | $t_1$ | $\cdots$ | $t_{14}$ | $t_{15}$ |
|---|---|---|---|---|
| $\downarrow$ | $\downarrow$ | | $\downarrow$ | $\downarrow$ |
| $w_0$ | $w_1$ | $\cdots$ | $w_{14}$ | $w_{15}$ |

| $t_0$ | $t_1$ | $\cdots$ | $t_{30}$ | $t_{31}$ |
|---|---|---|---|---|
| $\downarrow$ | $\downarrow$ | | $\downarrow$ | $\downarrow$ |
| $w_0$ | $w_1$ | $\cdots$ | $u_{14}$ | $u_{15}$ |

# Introduction
Implementation

The idea is to process integers at CUDA block-level, requiring:

- Sequentialization factor.

| $t_0$ | $t_1$ | $\cdots$ | | | | $t_0$ | | $\cdots$ | $t_{1023}$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ | | $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Coalesced transactions to global memory.

| $t_0$ | | $\cdots$ | $t_{1023}$ | | | $t_0$ | $t_1$ | $\cdots$ | $t_{1022}$ | $t_{1023}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ | | $w_0$ | $w_1$ | $\cdots$ | $w_{2046}$ | $w_{2047}$ |

- Multiple instances per block.

| $t_0$ | $t_1$ | $\cdots$ | $t_{14}$ | $t_{15}$ | | $t_0$ | $t_1$ | $\cdots$ | $t_{30}$ | $t_{31}$ |
|---|---|---|---|---|---|---|---|---|---|---|
| $w_0$ | $w_1$ | $\cdots$ | $w_{14}$ | $w_{15}$ | | $w_0$ | $w_1$ | $\cdots$ | $u_{14}$ | $u_{15}$ |

Introduction
oo

**Addition**
●ooooooooo

Multiplication
ooooooooooo

Division
oooooooooo

Conclusion
oo

# Addition

## Algorithm

### *badd* – bitwise-optimized addition of big integers

**Input:** Big integer $u$ and $v$ of size $m$ in base $B$
**Output:** Big integer $w$ of size $m$ in base $B$

```
1  (ws, cs) = map2 ⊕ us vs
2  pcs = scan_exc ⊗ 2 cs
3  w = map2 (λ w c → w + (c & 1) ) ws pcs
```

**where**

$$x \oplus y := (r, \ \mathrm{uint}(r < x) \mid (\mathrm{uint}(r == B - 1) \ll 1)), \quad \textbf{where } r = x + y \quad (2)$$

$$x \otimes y := (((x \ \& \ (y \gg 1)) \mid y) \ \& \ 1) \mid (x \ \& \ y \ \& \ 2) \quad (3)$$

## Addition
Implementation 1/2

Optimizations boils down to:

- Sequentialize the parallelism in excess.

- Segmented scan with flags integrated in the bitwise carry-overflow representation.

Introduction
00

**Addition**
0●00000000

Multiplication
00000000000

Division
0000000000

Conclusion
00

Addition
Implementation 1/2

Optimizations boils down to:

- Sequentialize the parallelism in excess.

- Segmented scan with flags integrated in the bitwise carry-overflow representation.

Introduction
○○

**Addition**
○○●○○○○○○○○

Multiplication
○○○○○○○○○○○

Division
○○○○○○○○○○

Conclusion
○○

# Addition
Implementation 2/2

Implementation revolves around a scan:

# Addition
## CUDA 1/3

Main addition kernel:

```
1   template<class B, uint32_t m, uint32_t q, uint32_t ipb>
2   __global__ void
3   baddKer3(typename B::uint_t* as, typename B::uint_t* bs,
4                                    typename B::uint_t* rs) {
5       using uint_t  = typename B::uint_t;
6       using carry_t = typename B::carry_t;
7       uint_t ass[q];
8       uint_t bss[q];
9       cpGlb2Reg<uint_t,m,q,ipb>(as, ass);
10      cpGlb2Reg<uint_t,m,q,ipb>(bs, bss);
11      __syncthreads();
12
13      uint_t rss[q];
14      __shared__ carry_t shmem[m*ipb];
15      baddKer3Run<B,m,q,ipb>(ass, bss, rss, shmem);   ⋆ badd
16      cpReg2Glb<uint_t,m,q,ipb>(rss, rs);
17  }
```

Introduction
oo

**Addition**
ooooooooooo

Multiplication
ooooooooooo

Division
ooooooooooo

Conclusion
oo

# Addition
## CUDA 2/3

Step 1.

```
18   uint_t css[q];
19   carry_t acc = threadIdx.x % (m/q) == 0
20       ? SegCarryProp<B>::setFlag(SegCarryProp<B>::identity())
21       : SegCarryProp<B>::identity();
22   #pragma unroll
23   for(int i=0; i<q; i++) {
24       rss[i] = ass[i] + bss[i];
25       css[i] = ((carry_t) (rss[i] < ass[i]))
26           | (((carry_t) (rss[i] == B::HIGHEST)) << 1);
27       acc = SegCarryProp<B>::apply(acc, css[i]);
28   }
29   shmem[threadIdx.x] = acc;
30   __syncthreads();
```

# Addition
## CUDA 3/3

Step 2.

```
31    acc = scanExcBlock< SegCarryProp<B> >(shmem, threadIdx.x);
32    acc = threadIdx.x % (m/q) == 0 ? SegCarryProp<B>::identity()
33                                   : acc;
```

Step 3.

```
34    #pragma unroll
35    for(int i=0; i<q; i++) {
36        rss[i] += (acc & 1);
37        acc = SegCarryProp<B>::apply(acc, css[i]);
38    }
```

## Addition
Futhark 1/3

Main addition function:

```
1  def baddV3 [ipb][m]
2  (us: [ipb*(4*m)]ui) (vs: [ipb*(4*m)]ui) : [ipb*(4*m)]ui =
3
4    let cp2sh (i: i64) = #[unsafe]
5      let str = ipb*m
6      in ((us[i], us[str + i], us[2*str + i], us[3*str + i])
7          ,(vs[i], vs[str + i], vs[2*str + i], vs[3*str + i]))
8
9    let (uss, vss) = map cp2sh (0..<ipb*m) |> unzip
10   let (u1s, u2s, u3s, u4s) = unzip4 uss
11   let (v1s, v2s, v3s, v4s) = unzip4 vss
12   let ush = u1s ++ u2s ++ u3s ++ u4s
13   let vsh = v1s ++ v2s ++ v3s ++ v4s
14
15   in baddV3Run ush vsh :> [ipb*(4*m)]ui    * badd
```

# Addition
### Futhark 2/3

Step 1.

```
16   let (ws, cs, accs) = unzip3 <| imap (0..<ipb*m) (\ i ->
17     let i4 = i*4
18     let (u1,u2,u3,u4) = (us[i4], us[i4+1], us[i4+2], us[i4+3])
19     let (v1,v2,v3,v4) = (vs[i4], vs[i4+1], vs[i4+2], vs[i4+3])
20     let (w1,w2,w3,w4) = (u1 + v1, u2 + v2, u3 + v3, u4 + v4)
21     let (c1,c2,c3,c4) = (carryAug w1 u1, carryAug w2 u2,
22                           carryAug w3 u3, carryAug w4 u4)
23     let c1 = (boolToCt (i % m == 0)) << 2 | c1
24     let acc = carryProp c1 <| carryProp c2 <| carryProp c3 c4
25     in ((w1, w2, w3, w4), (c1, c2, c3, c4), acc))
```

Step 2.

```
27   let pcs = scanExc carryPropSeg carryPropE accs
```

# Addition
## Futhark 3/3

Step 3.

```
28  let (wi1s, wi2s, wi3s, wi4s) = imap4 ws cs pcs (0..<ipb*m)
29    (\ (w1, w2, w3, w4) (c1, c2, c3, _) acc1 i ->
30        let acc1 = if i % m == 0 then carryPropE else acc1
31        let acc2 = carryProp acc1 c1
32        let acc3 = carryProp acc2 c2
33        let acc4 = carryProp acc3 c3
34        in ((w1 + fromCt (acc1 & 1), i*4),
35            (w2 + fromCt (acc2 & 1), i*4+1),
36            (w3 + fromCt (acc3 & 1), i*4+2),
37            (w4 + fromCt (acc4 & 1), i*4+3))) |> unzip
38
39  let (ws, inds) = unzip <| wi1s ++ wi2s ++ wi3s ++ wi4s
40  in scatter (replicate (ipb*(4*m)) 0) inds ws
```

## Addition
Evaluation

One addition and ten additions of base uint64_t in GB/s:

| Bits | Insts | CGBN1 | CUDA1 | Fut1 | CGBN10 | CUDA10 | Fut10 |
|------|-------|-------|-------|------|--------|--------|-------|
| $2^{18}$ | $2^{14}$ | 62 | 161 | – | 25 | 92 | – |
| $2^{17}$ | $2^{15}$ | 67 | 163 | – | 60 | 109 | – |
| $2^{16}$ | $2^{16}$ | 19 | 166 | 146 | 73 | 124 | 24 |
| $2^{15}$ | $2^{17}$ | 19 | 166 | 168 | 45 | 124 | 29 |
| $2^{14}$ | $2^{18}$ | 84 | 166 | 168 | 97 | 124 | 29 |
| $2^{13}$ | $2^{19}$ | 164 | 165 | 168 | 162 | 123 | 29 |
| $2^{12}$ | $2^{20}$ | 165 | 166 | 168 | 164 | 123 | 29 |
| $2^{11}$ | $2^{21}$ | 164 | 166 | 168 | 161 | 124 | 29 |
| $2^{10}$ | $2^{22}$ | 156 | 166 | 168 | 152 | 124 | 29 |
| $2^{9}$ | $2^{23}$ | 118 | 167 | 168 | 113 | 124 | 29 |

Introduction
○○

Addition
○○○○○○○○○○

Multiplication
●○○○○○○○○○○○

Division
○○○○○○○○○○

Conclusion
○○

# Classical Multiplication
## Algorithm 1/2

### Classical multiplication of big integers

Multiplying integer $u \in \mathbb{N}$ by $v \in \mathbb{N}$ in base $B$ and $m$ digits, is classically computed by:

$$u \cdot v = \sum_{k=0}^{m-1} \left( \sum_{\substack{0 \le i,j < m \\ i+j=k}} u_i \cdot v_j \right) B^k \qquad (4)$$

| $w_0$ | $w_1$ | $w_2$ | $\ldots$ | $w_{m-1}$ |
|-------|-------|-------|----------|-----------|
| $\overline{\overline{0}}$ | $\overline{\overline{c_0}}$ | $\overline{\overline{c_1}}$ | $\ldots$ | $\overline{\overline{c_{m-2}}}$ |
| $+$ | $+$ | $+$ | | $+$ |
| $u_0 v_0$ | $u_0 v_1$ | $u_0 v_2$ | $\ldots$ | $u_0 v_{m-1}$ |
| | $+$ | $+$ | | $+$ |
| | $u_1 v_0$ | $u_1 v_1$ | $\ldots$ | $u_1 v_{m-2}$ |
| | | $+$ | | $+$ |
| | | $u_2 v_0$ | $\ldots$ | $u_2 v_{m-3}$ |
| | | | $\ddots$ | $\vdots$ |
| | | | | $u_{m-1} v_0$ |

# Classical Multiplication
## Algorithm 2/2

### convmul – work-balanced classical multiplication by convolutions

Represent each convolution (column) by a *low*, *high*, and *carry* part.

Thread $k \in \{0, 1, \ldots, (m/2) - 1\}$ handles convolution $k_1 = k$ and $k_2 = m - 1 - k$.

The convolutions are added in the following pattern:

| $l_0^0$ | $l_1^1$ | $l_2^2$ | $l_3^3$ | $\cdots$ | $l_{m-4}^3$ | $l_{m-3}^2$ | $l_{m-2}^1$ | $l_{m-1}^0$ |
|---|---|---|---|---|---|---|---|---|
| + | + | + | + | | + | + | + | + |
| 0 | $h_0^0$ | $h_1^1$ | $h_2^2$ | $\cdots$ | $h_{m-5}^4$ | $h_{m-4}^3$ | $h_{m-3}^2$ | $h_{m-2}^1$ |
| + | + | + | + | | + | + | + | + |
| 0 | 0 | $c_0^0$ | $c_1^1$ | $\cdots$ | $c_{m-6}^5$ | $c_{m-5}^4$ | $c_{m-4}^3$ | $c_{m-3}^2$ |

# Classical Multiplication
Implementation

Optimizations boils down to:

- Double word-sizes to reduce work (for CUDA).
- Less communication but more sequentialization.

# Classical Multiplication
Implementation

Optimizations boils down to:

- Double word-sizes to reduce work (for CUDA).

- Less communication but more sequentialization.

Introduction
○○
Addition
○○○○○○○○○○
**Multiplication**
○○●○○○○○○○○○○
Division
○○○○○○○○○○
Conclusion
○○

# Classical Multiplication
## Implementation

Optimizations boils down to:

- Double word-sizes to reduce work (for CUDA).

- Less communication but more sequentialization.

# Classical Multiplication
## Implementation

Optimizations boils down to:

- Double word-sizes to reduce work (for CUDA).

- Less communication but more sequentialization.

# Classical Multiplication
## CUDA 1/3

The lower convolution:

```
1  ubig_t lh0[2]; lh0[0] = 0; lh0[1] = 0;
2  ubig_t lh1[2]; lh1[0] = 0; lh1[1] = 0;
3
4  int k1 = threadIdx.x*2;
5  int k1_start = (k1/m) * m;
6
7  for (int i=k1_start; i<=k1; i++) {
8      int j = k1 - i + k1_start;
9      uint_t a = shmem_as[i];
10     iterate<B>(a, shmem_bs[j], lh0);
11     iterate<B>(a, shmem_bs[j+1], lh1);
12  }
13  iterate<B>(shmem_as[k1+1], shmem_bs[k1_start], lh1);
14
15  combine2<B>(lh0, lh1, lhck1);
```

# Classical Multiplication
## CUDA 2/3

```
16   // iterate(uint_t a, uint_t b, ubig_t* lh)
17     ubig_t ab = ((ubig_t) a) * ((ubig_t) b);
18     lh[0] += ab & ((ubig_t) B::HIGHEST);
19     lh[1] += ab >> B::bits;
20
21   // combine2(ubig_t* lh0, ubig_t* lh1, uint_t* lhc)
22     uint_t h0t = (uint_t) lh0[1];
23     uint_t h0  = h0t + ((uint_t) (lh0[0] >> B::bits));
24     uint_t c0  = ((uint_t) (lh0[1] >> B::bits)) + (h0 < h0t);
25
26     uint_t h1t = (uint_t) lh1[1];
27     uint_t h1  = h1t + ((uint_t) (lh1[0] >> B::bits));
28     uint_t c1  = ((uint_t) (lh1[1] >> B::bits)) + (h1 < h1t);
29
30     lhc[0] = (uint_t) lh0[0];
31     lhc[1] = h0 + ((uint_t) lh1[0]);
32     lhc[2] = c0 + h1 + (lhc[1] < h0);
33     lhc[3] = c1 + (lhc[2] < h1);
```

Introduction
oo

Addition
oooooooooo

**Multiplication**
ooooooooooooo

Division
oooooooooo

Conclusion
oo

# Classical Multiplication
## CUDA 3/3

Memory layout:

# Classical Multiplication
## Futhark 1/3

Presented solution was sub-optimal.

Fixed by tagging parts with their index in the convolution function:

```
1  let s1 = (k1+2) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
2  let s2 = (k2+1) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
3  in
4  ((l1,  lh1,   l2,   lh2), (hc1,      cc1,      hc2,      cc2    ),
5   (k1,  k1+1, k2-1, k2), (s1|k1+2, s1|k1+3, s2|k2+1, s2|k2+2))
```

However, the operator is still slow for larger integers.

Inspired by [1], the shared memory is piped to opaque-function.

Significant speedup, but still slower than the basic version.

Introduction
oo

Addition
oooooooooo

Multiplication
ooooooo●oooo

Division
oooooooooo

Conclusion
oo

# Classical Multiplication
## Futhark 1/3

Presented solution was sub-optimal.

Fixed by tagging parts with their index in the convolution function:

```
1  let s1 = (k1+2) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
2  let s2 = (k2+1) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
3  in
4  ((l1, lh1, l2, lh2), (hc1, cc1, hc2, cc2),
5   (k1, k1+1, k2-1, k2), (s1|k1+2, s1|k1+3, s2|k2+1, s2|k2+2))
```

However, the operator is still slow for larger integers.

Inspired by [1], the shared memory is piped to opaque-function.

Significant speedup, but still slower than the basic version.

Introduction
oo

Addition
oooooooooo

**Multiplication**
ooooooo●oooo

Division
oooooooooo

Conclusion
oo

## Classical Multiplication
Futhark 1/3

Presented solution was sub-optimal.

Fixed by tagging parts with their index in the convolution function:

```
1  let  s1 = (k1+2) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
2  let  s2 = (k2+1) % (4*m) != 0 |> i64.bool |> (\i -> i - 1)
3  in
4  ((l1, lh1, l2, lh2), (hc1, cc1, hc2, cc2 ),
5   (k1, k1+1, k2-1, k2), (s1|k1+2, s1|k1+3, s2|k2+1, s2|k2+2))
```

However, the operator is still slow for larger integers.

Inspired by [1], the shared memory is piped to opaque-function.

Significant speedup, but still slower than the basic version.

Introduction
oo

Addition
oooooooooo

Multiplication
ooooooooo●ooo

Division
oooooooooo

Conclusion
oo

# Classical Multiplication
Futhark 2/3

```
1   let CONV (us: [] ui) (vs: [] ui) (tid: i64) = #[unsafe]
2     let k1 = tid
3     let k1_start = (k1 / (2*m)) * (2*m)
4     let lhc1 : (ui, ui, ui) =
5     loop (l, h, c) = (0, 0, 0)
6     for i < k1 + 1 − k1_start
7     do let j = k1 − i
8         let lr = us[i+k1_start] * vs[j]
9         let hr = mulHigh us[i+k1_start] vs[j]
10        let ln = l + lr
11        let hn = h + hr + (fromBool (ln < l))
12        let cn = c + (fromBool (hn < h))
13        in (ln, hn, cn)
14
15    let k2 = ipb*2*m−1 − k1
16    ...
17    in (lhc1, lhc2)
```

# Classical Multiplication
## Futhark 3/3

```
1   let ush = map (\i -> us[i]) (0..<ipb*(2*m))
2   let vsh = map (\i -> vs[i]) (0..<ipb*(2*m))
3
4   let (lhcs1, lhcs2) = map (CONV ush vsh) (0..<ipb*m) |> unzip
5   let (ls1, hs1, cs1) = unzip3 lhcs1
6   let (ls2, hs2, cs2) = unzip3 <| reverse lhcs2
7   let ls = ls1 ++ ls2 :> [ipb*(2*m)]ui
8   let hs = hs1 ++ hs2 :> [ipb*(2*m)]ui
9   let hs = map (\ i -> if i % (2*m) == 0 then 0 else hs[i-1] )
10               (0..<ipb*(2*m))
11  let cs = cs1 ++ cs2 :> [ipb*(2*m)]ui
12  let cs = map (\ i -> if i % (2*m) <= 1 then 0 else cs[i-2] )
13               (0..<ipb*(2*m))
14
15  in  baddV4 ls hs |> baddV4 cs
```

## Multiplication
Evaluation 1/2

One multiplication of base `uint64_t` in Gu32ops:

| Bits | Insts | CGBN | CUDA | FutOldQ2 | FutOldQ4 | FutNewQ4 | FutNewQ2 |
|------|-------|------|------|----------|----------|----------|----------|
| $2^{18}$ | $2^{14}$ | – | – | – | – | – | – |
| $2^{17}$ | $2^{15}$ | 1 | 1150 | – | – | – | – |
| $2^{16}$ | $2^{16}$ | 35 | 2039 | 974 | – | 1263 | 1453 |
| $2^{15}$ | $2^{17}$ | 116 | 3471 | 1674 | 482 | 2108 | 2423 |
| $2^{14}$ | $2^{18}$ | 217 | 5515 | 2671 | 693 | 3264 | 3697 |
| $2^{13}$ | $2^{19}$ | 340 | 8082 | 3880 | 984 | 4559 | 4947 |
| $2^{12}$ | $2^{20}$ | 526 | 10475 | 4931 | 1281 | 5467 | 5786 |
| $2^{11}$ | $2^{21}$ | 793 | 15745 | 3836 | 1899 | 6946 | 5990 |
| $2^{10}$ | $2^{22}$ | 822 | 16554 | 2352 | 2492 | 7830 | 6203 |
| $2^{9}$ | $2^{23}$ | 496 | 16888 | 1122 | 2798 | 8134 | 5646 |

## Multiplication
Evaluation 2/2

Six multiplications of base `uint64_t` in Gu32ops:

| Bits | Insts | CGBN | CUDA | FutOldQ2 | FutOldQ4 | FutNewQ4 | FutNewQ2 |
|------|-------|------|------|----------|----------|----------|----------|
| $2^{18}$ | $2^{14}$ | – | – | – | – | – | – |
| $2^{17}$ | $2^{15}$ | 11 | – | – | – | – | – |
| $2^{16}$ | $2^{16}$ | 888 | 1747 | 921 | – | – | – |
| $2^{15}$ | $2^{17}$ | 2832 | 2602 | 1595 | 350 | – | – |
| $2^{14}$ | $2^{18}$ | 4960 | 2696 | 1656 | 513 | 1211 | 1609 |
| $2^{13}$ | $2^{19}$ | 8625 | 4961 | 1872 | 778 | 1264 | 2194 |
| $2^{12}$ | $2^{20}$ | 13924 | 8981 | 3307 | 1029 | 1507 | 2616 |
| $2^{11}$ | $2^{21}$ | 23424 | 13717 | 3028 | 1505 | 2068 | 2681 |
| $2^{10}$ | $2^{22}$ | 37500 | 17513 | 2180 | 1946 | 2452 | 2859 |
| $2^{9}$ | $2^{23}$ | 70093 | 17079 | 1225 | 2156 | 2626 | 2678 |

# Division
## Algorithm 1/3

The intuition behind the division algorithm:

- Multiply the dividend with the of inverse divisor.

- Use shifts to represent the inverse as a big integer.

- Approximate the shifted inverse by Newton's Method using [2].

- Compute quotient and remainder to adjust approximation.

### Quotient of big integers by shifted inverse

We define the quotient of big integers $u \leq B^{h \in \mathbb{N}}$ and $v$ in base $B$ using [2] as:

$$u \text{ quo } v = \text{shift}_{-h} (u \cdot \text{shinv}_h v) + \delta, \quad \textbf{where } \delta \in \{0, 1\} \qquad (5)$$

$$\text{shift}_{n \in \mathbb{Z}} u = \lfloor u \cdot B^n \rfloor \qquad \qquad \text{shinv}_{n \in \mathbb{N}} v = \lfloor B^n / v \rfloor \qquad (6)$$

# Division
Algorithm 1/3

The intuition behind the division algorithm:

- Multiply the dividend with the of inverse divisor.

- Use shifts to represent the inverse as a big integer.

- Approximate the shifted inverse by Newton's Method using [2].

- Compute quotient and remainder to adjust approximation.

## Quotient of big integers by shifted inverse

We define the quotient of big integers $u \leq B^{h \in \mathbb{N}}$ and $v$ in base $B$ using [2] as:

$$u \text{ quo } v = \text{shift}_{-h}\left(u \cdot \text{shinv}_h v\right) + \delta, \quad \text{where } \delta \in \{0, 1\} \tag{5}$$

$$\text{shift}_{n \in \mathbb{Z}} u = \lfloor u \cdot B^n \rfloor \qquad \text{shinv}_{n \in \mathbb{N}} v = \lfloor B^n / v \rfloor \tag{6}$$

# Division
## Algorithm 1/3

The intuition behind the division algorithm:

- Multiply the dividend with the of inverse divisor.

- Use shifts to represent the inverse as a big integer.

- Approximate the shifted inverse by Newton's Method using [2].

- Compute quotient and remainder to adjust approximation.

### Quotient of big integers by shifted inverse

We define the quotient of big integers $u \leq B^{h \in \mathbb{N}}$ and $v$ in base $B$ using [2] as:

$$u \text{ quo } v = \text{shift}_{-h}(u \cdot \text{shinv}_h v) + \delta, \quad \text{where } \delta \in \{0, 1\} \tag{5}$$

$$\text{shift}_{n \in \mathbb{Z}} u = \lfloor u \cdot B^n \rfloor \qquad \text{shinv}_{n \in \mathbb{N}} v = \lfloor B^n / v \rfloor \tag{6}$$

# Division
## Algorithm 1/3

The intuition behind the division algorithm:

- Multiply the dividend with the of inverse divisor.

- Use shifts to represent the inverse as a big integer.

- Approximate the shifted inverse by Newton's Method using [2].

- Compute quotient and remainder to adjust approximation.

### Quotient of big integers by shifted inverse

We define the quotient of big integers $u \leq B^{h \in \mathbb{N}}$ and $v$ in base $B$ using [2] as:

$$u \text{ quo } v = \text{shift}_{-h} (u \cdot \text{shinv}_h v) + \delta, \quad \text{where } \delta \in \{0, 1\} \tag{5}$$

$$\text{shift}_{n \in \mathbb{Z}} \ u = \lfloor u \cdot B^n \rfloor \qquad \text{shinv}_{n \in \mathbb{N}} \ v = \lfloor B^n/v \rfloor \tag{6}$$

# Division
## Algorithm 1/3

The intuition behind the division algorithm:

- Multiply the dividend with the of inverse divisor.

- Use shifts to represent the inverse as a big integer.

- Approximate the shifted inverse by Newton's Method using [2].

- Compute quotient and remainder to adjust approximation.

### Quotient of big integers by shifted inverse

We define the quotient of big integers $u \leq B^{h \in \mathbb{N}}$ and $v$ in base $B$ using [2] as:

$$u \text{ quo } v = \text{shift}_{-h} (u \cdot \text{shinv}_h v) + \delta, \quad \textbf{where } \delta \in \{0, 1\} \tag{5}$$

$$\text{shift}_{n \in \mathbb{Z}} u = \lfloor u \cdot B^n \rfloor \qquad \text{shinv}_{n \in \mathbb{N}} v = \lfloor B^n / v \rfloor \tag{6}$$

# Division
Algorithm 2/3

The algorithm of [2] is based on the Netwon iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} = x_i + \left(x_i - \frac{v}{u}x_i^2\right), \quad \text{where } x \in \mathbb{R} \text{ and } f(x) = \frac{u}{x} - v \qquad (7)$$

It is modified in [2] w.r.t. three aspects:

- It is discretized to integers, so we compute $w \in \mathbb{Z}$ rather than $x \in \mathbb{R}$.

- $u$ is specialized to $B^h$ (where $B$ is the base and $u \leq B^h$).

- Use a shift rather than a division since $v/(B^h) = v \cdot B^{-h} \geq \text{shift}_{-h} \, v$.

The Newton iteration becomes:

$$w_{i+1} = w_i + \text{shift}_{-h}\left(\text{shift}_h \, w_i - v \cdot w_i^2\right) = w_i + w_i\lfloor B^h - v \cdot w_i \rfloor B^{-h} \qquad (8)$$

# Division
Algorithm 2/3

The algorithm of [2] is based on the Netwon iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} = x_i + \left(x_i - \frac{v}{u}x_i^2\right), \quad \text{where } x \in \mathbb{R} \text{ and } f(x) = \frac{u}{x} - v \qquad (7)$$

It is modified in [2] w.r.t. three aspects:

- It is discretized to integers, so we compute $w \in \mathbb{Z}$ rather than $x \in \mathbb{R}$.
- $u$ is specialized to $B^h$ (where $B$ is the base and $u \leq B^h$).
- Use a shift rather than a division since $v/(B^h) = v \cdot B^{-h} \geq \text{shift}_{-h}\ v$.

The Newton iteration becomes:

$$w_{i+1} = w_i + shift_{-h}\left(shift_h\ w_i - v \cdot w_i^2\right) = w_i + w_i\lfloor B^h - v \cdot w_i\rfloor B^{-h} \qquad (8)$$

# Division
Algorithm 2/3

The algorithm of [2] is based on the Netwon iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} = x_i + \left(x_i - \frac{v}{u}x_i^2\right), \quad \text{where } x \in \mathbb{R} \text{ and } f(x) = \frac{u}{x} - v \qquad (7)$$

It is modified in [2] w.r.t. three aspects:

- It is discretized to integers, so we compute $w \in \mathbb{Z}$ rather than $x \in \mathbb{R}$.

- $u$ is specialized to $B^h$ (where $B$ is the base and $u \leq B^h$).

- Use a shift rather than a division since $v/(B^h) = v \cdot B^{-h} \geq \text{shift}_{-h} \, v$.

The Newton iteration becomes:

$$w_{i+1} = w_i + \text{shift}_{-h}\left(\text{shift}_h \, w_i - v \cdot w_i^2\right) = w_i + w_i\lfloor B^h - v \cdot w_i\rfloor B^{-h} \qquad (8)$$

# Division
## Algorithm 2/3

The algorithm of [2] is based on the Netwon iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} = x_i + \left( x_i - \frac{v}{u}x_i^2 \right), \quad \text{where } x \in \mathbb{R} \text{ and } f(x) = \frac{u}{x} - v \qquad (7)$$

It is modified in [2] w.r.t. three aspects:

- It is discretized to integers, so we compute $w \in \mathbb{Z}$ rather than $x \in \mathbb{R}$.

- $u$ is specialized to $B^h$ (where $B$ is the base and $u \leq B^h$).

- Use a shift rather than a division since $v/(B^h) = v \cdot B^{-h} \geq \text{shift}_{-h} \, v$.

The Newton iteration becomes:

$$w_{i+1} = w_i + shift_{-h} \left( shift_h \, w_i - v \cdot w_i^2 \right) = w_i + w_i \lfloor B^h - v \cdot w_i \rfloor B^{-h} \qquad (8)$$

# Division
## Algorithm 2/3

The algorithm of [2] is based on the Netwon iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} = x_i + \left(x_i - \frac{v}{u}x_i^2\right), \quad \text{where } x \in \mathbb{R} \text{ and } f(x) = \frac{u}{x} - v \qquad (7)$$

It is modified in [2] w.r.t. three aspects:

- It is discretized to integers, so we compute $w \in \mathbb{Z}$ rather than $x \in \mathbb{R}$.

- $u$ is specialized to $B^h$ (where $B$ is the base and $u \leq B^h$).

- Use a shift rather than a division since $v/(B^h) = v \cdot B^{-h} \geq \mathtt{shift}_{-h}\ v$.

The Newton iteration becomes:

$$w_{i+1} = w_i + shift_{-h}\left(shift_h\ w_i - v \cdot w_i^2\right) = w_i + w_i\lfloor B^h - v \cdot w_i\rfloor B^{-h} \qquad (8)$$

## Division
Algorithm 3/3

Instead of showing the algorithm for computing the division and the shifted inverse, let us look at the Futhark implementation.

In the thesis, the Futhark implementation was invalid.

Now it validates!

(Without the divisor prefixes and shorter iterates optimizations.)

Hence all arithmetics are in full length.

The analysis in [2] gives work $O(\log(h - k)(M(h) + M(|h/2 - k|)))$.

If we assume $h = m + k$, we get $O(\log(m)M(m))$.

# Division
## Algorithm 3/3

Instead of showing the algorithm for computing the division and the shifted inverse, let us look at the Futhark implementation.

In the thesis, the Futhark implementation was invalid.

Now it validates!

(Without the divisor prefixes and shorter iterates optimizations.)

Hence all arithmetics are in full length.

The analysis in [2] gives work $O(\log(h - k)(M(h) + M(|h/2 - k|)))$.

If we assume $h = m + k$, we get $O(\log(m)M(m))$.

# Division
Algorithm 3/3

Instead of showing the algorithm for computing the division and the shifted inverse, let us look at the Futhark implementation.

In the thesis, the Futhark implementation was invalid.

Now it validates!

(Without the divisor prefixes and shorter iterates optimizations.)

Hence all arithmetics are in full length.

The analysis in [2] gives work $O(\log(h-k)(M(h)+M(|h/2-k|)))$.

If we assume $h = m + k$, we get $O(\log(m)M(m))$.

## Division
Implementation 1/5

Revisions to section 9.2 and section 9.3 of the thesis:

- 2 guard digits are sufficient.

- The thesis states that $v > B^h$ corresponds to:

$$\exists i \in \mathbb{N}. \ (h < i < m \land v[i] \neq 0) \lor (h = i < m \land v[i] > 1) \qquad (9)$$

This is not true. E.g. it fails on $[1, 0, 0, 1] > B^3$.

Instead, define $v > B^h$ as $\neg(v < B^h \lor v = B^h)$.

Define $v < B^h$ as $\forall i \in \mathbb{N}. \ i \geq h \lor v[i] = 0$.

Define $v = B^h$ as $\forall i \in \mathbb{N}. \ (i = h \land v[i] = 1) \lor (i \neq h \land v[i] = 0)$.

Introduction
oo

Addition
oooooooooo

Multiplication
ooooooooooo

**Division**
oooooooooo

Conclusion
oo

## Division
Implementation 1/5

Revisions to section 9.2 and section 9.3 of the thesis:

- 2 guard digits are sufficient.

- The thesis states that $v > B^h$ corresponds to:

$$\exists i \in \mathbb{N}. \ (h < i < m \land v[i] \neq 0) \lor (h = i < m \land v[i] > 1) \qquad (9)$$

This is not true. E.g. it fails on $[1, 0, 0, 1] > B^3$.

Instead, define $v > B^h$ as $\neg(v < B^h \lor v = B^h)$.

Define $v < B^h$ as $\forall i \in \mathbb{N}. \ i \geq h \lor v[i] = 0$.

Define $v = B^h$ as $\forall i \in \mathbb{N}. \ (i = h \land v[i] = 1) \lor (i \neq h \land v[i] = 0)$.

## Division
Implementation 1/5

Revisions to section 9.2 and section 9.3 of the thesis:

- 2 guard digits are sufficient.
- The thesis states that $v > B^h$ corresponds to:

$$\exists i \in \mathbb{N}. \ (h < i < m \land v[i] \neq 0) \lor (h = i < m \land v[i] > 1) \qquad (9)$$

  This is not true. E.g. it fails on $[1, 0, 0, 1] > B^3$.

  Instead, define $v > B^h$ as $\neg(v < B^h \lor v = B^h)$.

  Define $v < B^h$ as $\forall i \in \mathbb{N}. \ i \geq h \lor v[i] = 0$.

  Define $v = B^h$ as $\forall i \in \mathbb{N}. \ (i = h \land v[i] = 1) \lor (i \neq h \land v[i] = 0)$.

Introduction
○○

Addition
○○○○○○○○○○

Multiplication
○○○○○○○○○○○○○

**Division**
○○○○○●○○○○○

Conclusion
○○

# Division
Implementation 2/5

```
1  def div [m] (u: [m] ui) (v: [m] ui) : ([m] ui, [m] ui) =
2    let h = findh u -- u ≤ B^h
3    let k = findk v -- B^k ≤ v < B^(k+1)
4
5    let p = 2*(m + (i64.bool (k <= 1)) + (i64.bool (k == 0)))
6    let up = map (\ i -> if i < m then u[i] else 0 ) (iota p)
7    let vp = map (\ i -> if i < m then v[i] else 0 ) (iota p)
8
9    let (h, k, up, vp) =
10            if k == 1 then (h+1, k+1, shift 1 up, shift 1 vp)
11      else if k == 0 then (h+2, k+2, shift 2 up, shift 2 vp)
12      else              (h,   k,   up,          vp)
13
14    let q = shinv k vp h |> mul up |> shift (−h) |> take m
15    let r = mul v q |> sub u |> fst
16    in if not (lt r v)
17      then (add q (singleton m 1), fst (sub r v)) else (q, r)
```

# Division
Implementation 3/5

```
18  def shinv [m] (k: i64) (v: [m] ui) (h: i64) : [m] ui =
19    assert (k > 1) (
20          if gtBpow v h          then new m
21      else if gtBpow (muld v 2) h then singleton m 1
22      else if eqBpow v k          then bpow m (h - k)
23      else
24        let V = (toQi v[k-2])
25              + (toQi v[k-1] << (i64ToQi bits))
26              + (toQi v[k]   << (i64ToQi (2*bits)))
27        let W = ((0 - V) / V) + 1
28        let w = map (\i -> if i <= 1
29                           then fromQi (W >> (i64ToQi (bits*i)))
30                           else 0) (iota m)
31
32        in if h - k <= 2 then shift (h - k - 2) w
33                          else refine v w h k       )
```

Introduction
○○

Addition
○○○○○○○○○○

Multiplication
○○○○○○○○○○○○

**Division**
○○○○○○○●○○○

Conclusion
○○

# Division
Implementation 4/5

```
34   def refine [m] (v:[m]ui) (w:[m]ui) (h:i64) (k:i64) : [m]ui =
35     let g = 1
36     let h = h + g
37     let (w, _) =
38       loop (w, l) = (shift (h−k−2) w, 2) while h − k > l do
39       let w = step h v w 0 l 0
40       let l = i64.min (2*l−1) (h−k)
41       in (w, l)
42     in shift (−g) w
```

```
43   def step [m] (h: i64) (v: [m]ui) (w: [m]ui)
44               (n: i64) (l: i64)   (g:i64)   : [m]ui =
45     let (pwd, sign) = powdiff v w (h−n) (l−g)
46     let wpwdS = shift (2*n − h) (mul w pwd)
47     let wS = shift n w
48     in if sign then fst (sub wS wpwdS) else add wS wpwdS
```

# Division
Implementation 5/5

```
49   def powdiff [m] (v: [m]ui) (w: [m]ui)
50                   (h: i64)    (l: i64)   : ([m]ui, bool) =
51     let L = (prec v) + (prec w) − l + 1
52     in if (ez v) || (ez w) then (bpow m h, false)
53        else if L >= h then sub (bpow m h) (mul v w)
54        else let P = multmod v w L
55             in if ez P then (P, false)
56                else if P[L−1] == 0 then (P, true)
57                else sub (bpow m L) P
```

```
58   def multmod [m] (v: [m]ui) (w: [m]ui) (e: i64) : [m]ui =
59     let vw = mul (take e v) (take e w)
60     in map (\ i −> if i < e then vw[i] else 0 ) (iota m)
```

# Division
Evaluation 1/2

The implementation is not efficient:

- Batch processing results in error:

  "Known compiler limitation encountered. Sorry."

  Can be circumvented with attribute #[sequential_outer].

- It succeeds in generating intra-block version when run with:

  #[incremental_flattening(only_intra)]

  It runs significantly slower for sizes greater than 256 digits.

# Division
Evaluation 1/2

The implementation is not efficient:

- Batch processing results in error:

  "Known compiler limitation encountered. Sorry."

  Can be circumvented with attribute #[sequential_outer].

- It succeeds in generating intra-block version when run with:

  #[incremental_flattening(only_intra)]

  It runs significantly slower for sizes greater than 256 digits.

## Division
Evaluation 2/2

Furthermore, runtimes depend on precision rather than size.

Thus, incomparable to the evaluation method for other arithmetics.

However, the implementation has no problems compiling to C code.

Hence, we could use multicore backend and compare to GMP.

The difference in runtimes are so big that results are meaningless.

Conclusion: It is inefficient and not comparable to GMP or CGBN.

# Division
Evaluation 2/2

Furthermore, runtimes depend on precision rather than size.

Thus, incomparable to the evaluation method for other arithmetics.

However, the implementation has no problems compiling to C code.

Hence, we could use `multicore` backend and compare to GMP.

The difference in runtimes are so big that results are meaningless.

Conclusion: It is inefficient and not comparable to GMP or CGBN.

# Division
Evaluation 2/2

Furthermore, runtimes depend on precision rather than size.

Thus, incomparable to the evaluation method for other arithmetics.

However, the implementation has no problems compiling to C code.

Hence, we could use `multicore` backend and compare to GMP.

The difference in runtimes are so big that results are meaningless.

Conclusion: It is inefficient and not comparable to GMP or CGBN.

## Conclusion

This thesis has shown:

- How to compute exact addition, classical multiplication, and division of big integers in parallel.

- How to implement the arithmetic operators in Futhark and CUDA C++.

- Efficient implementations of addition and multiplication.

- Correct but inefficient implementation of division, while outlining the required steps to efficiency.

# Conclusion

This thesis has shown:

- How to compute exact addition, classical multiplication, and division of big integers in parallel.

- How to implement the arithmetic operators in Futhark and CUDA C++.

- Efficient implementations of addition and multiplication.

- Correct but inefficient implementation of division, while outlining the required steps to efficiency.

# Conclusion

This thesis has shown:

- How to compute exact addition, classical multiplication, and division of big integers in parallel.

- How to implement the arithmetic operators in Futhark and CUDA C++.

- Efficient implementations of addition and multiplication.

- Correct but inefficient implementation of division, while outlining the required steps to efficiency.

## Conclusion

This thesis has shown:

- How to compute exact addition, classical multiplication, and division of big integers in parallel.

- How to implement the arithmetic operators in Futhark and CUDA C++.

- Efficient implementations of addition and multiplication.

- Correct but inefficient implementation of division, while outlining the required steps to efficiency.

# Conclusion
References

📄 Cosmin E. Oancea and Stephen M. Watt.

GPU Implementations for Midsize Integer Addition and Multiplication, 2024.

arXiv:2405.14642.

📄 Stephen M. Watt.

Efficient Generic Quotients Using Exact Arithmetic.

2023.

arXiv:2304.01753.