

Základní paralelní operace

Tomáš Oberhuber

`tomas.oberhuber@fjfi.cvut.cz`

15. dubna 2024

Videa na Youtube:

Redukce

Prefix sum

Základní paralelní operace

Guy E. Blelloch, Vector Models for Data-Parallel Computing,
Cambridge: MIT press, 1990.

Co je paralelní redukce?

Definition

Mějme pole prvků a_1, \dots, a_n určitého typu a asociativní operaci \oplus . **Paralelní redukce** je paralelní aplikace asociativní operace \oplus na všechny prvky vstupního pole tak, že výsledkem je jeden prvek a stejného typu, pro který platí

$$a = a_1 \oplus a_2 \dots a_n.$$

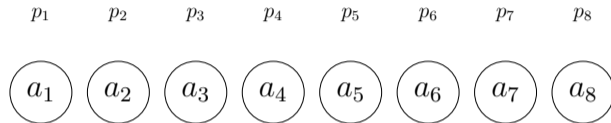
- ▶ porovnání polí, vektorů a řetězců
- ▶ skalární součin dvou vektorů
- ▶ l_p normy
- ▶ minimalální maximalní hodnota
- ▶ suma všech prvků pole

Co je paralelní redukce?

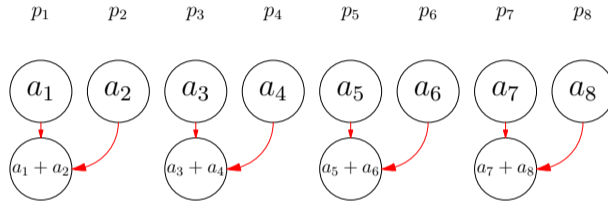
Příklad: Pro jednoduchost budeme uvažovat operaci sčítání. Pak jde o paralelní provedení tohoto kódu:

```
a = a[ 1 ];  
for( int i = 2; i <= n; i ++ )  
    a += a[ i ];
```

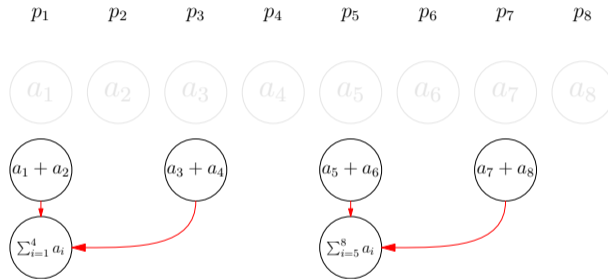
Paralelní provedení



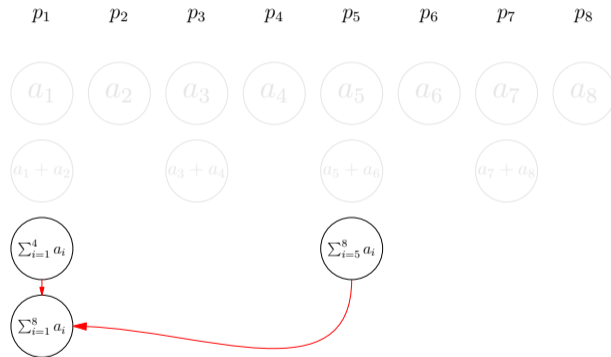
Paralelní provedení



Paralelní provedení



Paralelní provedení



Analýza paralelní provedení

Za předpokladu $n = p$, snadno vidíme, že platí:

- ▶ $T_S(n) = \theta(n)$
- ▶ $T_P(n) = \theta(\log n)$
- ▶ $S(n) = \theta\left(\frac{n}{\log n}\right) = O(n) = O(p)$
- ▶ $E(n) = \theta\left(\frac{1}{\log n}\right)$
 - ▶ pro velká n efektivita klesá k nule
- ▶ $C(n) = \theta(n \log n) = \omega(T_S(n))$
 - ▶ algoritmus není nákladově optimální

Provedeme úpravu algoritmu tak, aby byl nákladově optimální.

Nákladově optimální redukce

p_1



p_2



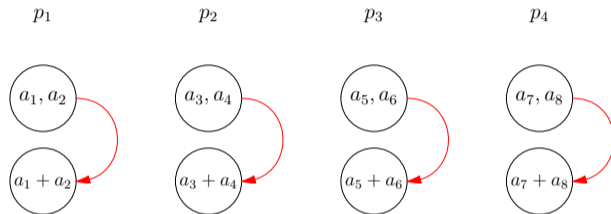
p_3



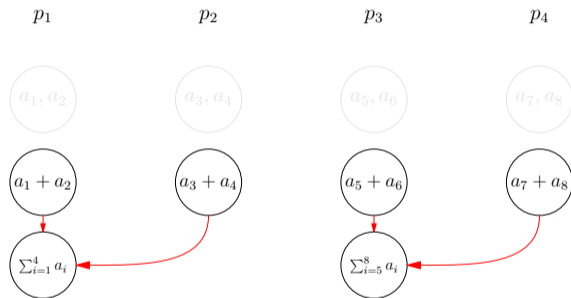
p_4



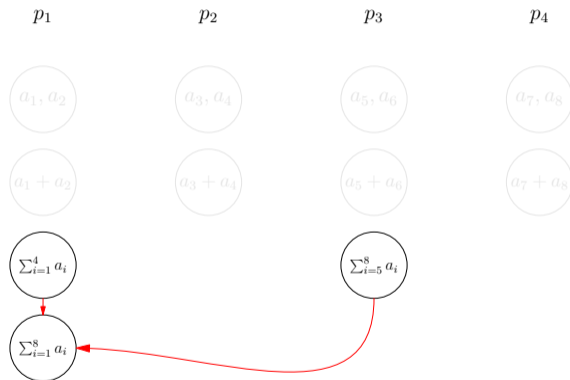
Nákladově optimální redukce



Nákladově optimální redukce



Nákladově optimální redukce



Nákladově optimální redukce

- ▶ nyní máme méně procesorů než prvků, tj. $p < n$.
- ▶ nejprve se provede sekvenční redukce $\frac{n}{p}$ prvků na každém procesoru $\rightarrow \theta\left(\frac{n}{p}\right)$
- ▶ a pak paralelní redukce na p procesorech $\rightarrow \theta(\log p)$
- ▶ celkem tedy $T_P(n, p) = \theta\left(\frac{n}{p} + \log p\right)$
- ▶ $S(n, p) = \frac{T_S(n)}{T_P(n, p)} = \theta\left(\frac{n}{\frac{n}{p} + \log p}\right) = \theta\left(\frac{p}{1 + \frac{p \log p}{n}}\right)$
- ▶ $E(n, p) = \frac{S(n, p)}{p} = \theta\left(\frac{1}{1 + \frac{p \log p}{n}}\right)$

Nákladově optimální redukce

Pak platí

- ▶ $C(n, p) = pT_P(n, p) = \theta(n + p \log p)$
- ▶ protože je $n = \Omega(p \log p)$, máme

$$C(n, p) = \theta(n) = \theta(T_S(n)).$$

Algoritmus je tedy nákladově optimální.

- ▶ toho jsme dosáhli na úkor počtu procesorů
- ▶ některé algoritmy nedokáží optimálně využít velký počet procesorů

Paralelní redukce na architekturách se sdílenou pamětí

- ▶ OpenMP má podporu pro redukci s základními operacemi –
`reduction (sum: +)`
- ▶ pokud chceme dělat redukci s jinou operací (`min`, `max`) nebo jiným, než základním typem (matice), je nutné ji provést explicitně
- ▶ tyto architektury mají většinou maximálně několik desítek procesorů
- ▶ redukci lze často provádět sekvenčně
- ▶ mezivýsledky se ukládají do nesdílených proměnných
 - ▶ pozor na cache coherence a false sharing - tyto proměnné by neměly být alokované v jednom bloku
- ▶ nakonec se zapíše do sdíleného pole, na kterém se provede redukce třeba i sekvenčně nultým procesem

Paralelní redukce na architekturách s distribuovanou pamětí

- ▶ standard MPI má velmi dobrou podporu pro redukci
- ▶ podporuje více operací a to i pro odvozené typy a struktury

Paralelní redukce na GPU v CUDA

M. Harris, *Optimizing Parallel Reduction in CUDA*, Nvidia.

Paralelní redukce v TNL

Jako příklad si vezmeme skalární součin:

```
1 float result( 0.0 );
2 for( int i = 0; i < size; i++ )
3     result += a[ i ] * b[ i ];
```

Nyní ho přepíšeme pomocí lambda funkcí:

```
1 auto fetch = [=] __cuda_callable__ (int i)->float { return a[i]*b[i]; };
2 auto reduction = [] __cuda_callable__ (float x, float y) -> float {
3     return x+y; };
4
5 float result( 0.0 );
6 for( int i = 0; i < size; i++ )
7     reduction = reduction( result, fetch( i ) );
```

Paralelní redukce v TNL

V TNL, pak for cyklus nahradíme funkcí `reduction`:

```
1  auto a_view = a.getView();
2  auto b_view = b.getView();
3  auto fetch = [=] __cuda_callable__ ( int i)->float {
4      return a_view[ i ] * b_view[ i ]; };
5  auto reduction = [] __cuda_callable__ (float x, float y) -> float {
6      return x + y; };
7
8  result = TNL::Algorithms::reduce< Device >( 0,
9                                              a.getSize(),
10                                             fetch,
11                                             reduction,
12                                             0.0 );
```

Poslední parametr je identický element pro danou operaci.

Paralelní redukce v TNL

Porovnání dvou vektorů pak může být provedeno takto:

```
1  auto a_view = a.getView();
2  auto b_view = b.getView();
3  auto fetch = [=] __cuda_callable__ (int i)->bool {
4      return a_view[ i ] == b_view[ i ]; };
5  auto reduction = [] __cuda_callable__ (float x, float y)->float {
6      return x && y; };
7
8  result = TNL::Algorithms::reduce< Device >( 0,
9                                              a.getSize(),
10                                             fetch,
11                                             reduction,
12                                             true );
```

Paralelní redukce v TNL

Největší prvek v absolutní hodnotě lze najít takto:

```
1 auto a_view = a.getView();
2 auto fetch = [=] __cuda_callable__ (int i)->float {
3     return abs( a_view[ i ] ); };
4 auto reduction = [] __cuda_callable__ (float x, float y)->float {
5     return max( x, y ); };
6
7 result = TNL::Algorithms::reduce< Device >( 0,
8                                             a.getSize(),
9                                             fetch,
10                                            reduction,
11                                            std::numeric_limits< float >::lowest() );
```

Nebo kompaktněji:

```
1 auto a_view = a.getView();
2 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(),
3     [=] __cuda_callable__ (int i)->float { return abs( a_view[ i ] ); },
4     TNL::Max() );
```


Paralelní redukce v TNL

Nalezení největšího prvku v absolutní hodnotě:

```
1  template< typename Device >
2  std::pair< double, int >
3  maximumNorm( const Vector< double, Device >& v )
4  {
5      auto view = v.getConstView();
6
7      auto fetch = [=] __cuda_callable__ ( int i ) {
8          return abs( view[ i ] ); };
9      auto reduction = [] __cuda_callable__ ( double& a,
10                                             const double& b,
11                                             int& aIdx,
12                                             const int& bIdx ) {
13          if( a < b ) {
14              a = b;
15              aIdx = bIdx;
16          }
17          else if( a == b && bIdx < aIdx )
18              aIdx = bIdx;
19      };
20      return reduceWithArgument< Device >( 0, view.getSize(),
21                                           fetch, reduction,
22                                           std::numeric_limits< double >::max() );
23 }
```

Paralelní redukce v TNL

Volání předchozí funkce vypadá takto:

```
1  int main( int argc, char* argv[] )
2  {
3      Vector< double, Devices::Host > host_v( 10 );
4      host_v.forAllElements( [] __cuda_callable__ ( int i, double& value ) {
5          value = i - 7; } );
6      std::cout << "host_v = " << host_v << std::endl;
7      auto maxNormHost = maximumNorm( host_v );
8      std::cout << "The maximum norm of the host vector elements is "
9          << maxNormHost.first
10         << " at position "
11         << maxNormHost.second << "." << std::endl;
12
13     #ifdef __CUDACC__
14         Vector< double, Devices::Cuda > cuda_v( 10 );
15         cuda_v.forAllElements( [] __cuda_callable__ ( int i, double& value ) {
16             value = i - 7; } );
17         std::cout << "cuda_v = " << cuda_v << std::endl;
18         auto maxNormCuda = maximumNorm( cuda_v );
19         std::cout << "The maximum norm of the device vector elements is "
20             << maxNormCuda.first
21             << " at position "
22             << maxNormCuda.second << "." << std::endl;
23     #endif
24     return EXIT_SUCCESS;
25 }
```

Paralelní redukce v TNL

Také je možné spojit dvě operace do jedné:

- ▶ přičtení vektoru
- ▶ výpočet normy tohoto vektoru.

Toto se často vyskytuje v Rungových-Kuttových metodách:

```
1 auto a_view = a.getView();
2 auto b_view = b.getView();
3 result = TNL::Algorithms::reduce< Device >( 0, a.getSize(),
4     [=] __cuda_callable__ (int i)->float mutable {
5         // mutable because we change a_view
6         float update = 0.5 * ( a_view[ i ] + b_view[ i ] );
7         a_view[ i ] += update;
8         return abs( update ); },
9     TNL::Plus() );
```

Paralelní redukce v Thrust

```
1  #include <thrust/reduce.h>
2  #include <thrust/device_vector.h>
3  #include <iostream>
4
5  int main() {
6      // Example data
7      thrust::device_vector<int> data{1, 2, 3, 4, 5};
8
9      // Perform reduction
10     int sum = thrust::reduce(data.begin(),
11                             data.end(),
12                             0,
13                             thrust::plus<int>());
14
15     std::cout << "Sum: " << sum << std::endl;
16
17     return 0;
18 }
```

Paralelní redukce v Kokkos

```
1  #include <Kokkos_Core.hpp>
2  #include <iostream>
3
4  int main(int argc, char* argv[]) {
5      Kokkos::initialize(argc, argv);
6      {
7          int N = 5;
8          Kokkos::View<int*> data("data", N);
9          Kokkos::parallel_for("InitData", N, KOKKOS_LAMBDA(const int i) {
10             data(i) = i + 1;
11         });
12
13         int sum = 0;
14         // Use a lambda for the reduction
15         Kokkos::parallel_reduce("SumReduce", N,
16             KOKKOS_LAMBDA(const int i, int& lsum) {
17                 lsum += data(i);
18             },
19             sum);
20
21         std::cout << "Sum: " << sum << std::endl;
22     }
23     Kokkos::finalize();
24
25     return 0;
26 }
```


Paralelní redukce a MapReduce v Hadoop

*"Hadoop je framework obsahující sadu opensource softwarových komponent určených pro zpracování velkého množství nestrukturovaných a **distribuovaných** dat v řádech petabytů a exabytů."*¹

Obsahuje tyto nástroje:

- ▶ **Hadoop Distributed File System (HDFS)** - distribuovaný filesystem.
- ▶ **MapReduce** - programovací model pro paralelní zpracování velkých datových sad na distribuovaných klastrech.
- ▶ **YARN (Yet Another Resource Negotiator)** - distribuovaná plánovač úloh.

Ukážeme se příklad na MapReduce v Javě.

¹Zdroj Wikipedie.

Paralelní redukce a MapReduce v Hadoop

Nejprve implementujeme `Mapper` \equiv obdoba funkce `fetch` v TNL.

```
1  import org.apache.hadoop.io.IntWritable;
2  import org.apache.hadoop.io.LongWritable;
3  import org.apache.hadoop.io.Text;
4  import org.apache.hadoop.mapreduce.Mapper;
5
6  import java.io.IOException;
7
8  public class WordCountMapper
9      extends Mapper<LongWritable, Text, Text, IntWritable> {
10
11     private final static IntWritable one = new IntWritable(1);
12     private Text word = new Text();
13
14     public void map(LongWritable key, Text value, Context context)
15         throws IOException, InterruptedException {
16         String[] words = value.toString().split("\\s+");
17         for (String str : words) {
18             word.set(str);
19             context.write(word, one);
20         }
21     }
22 }
```

- `Mapper` čte vstupní text `value`, dělí ho na slova a na každé slovo `word` vrací dvojici klíč-hodnota ve tvaru `word-1`.

Paralelní redukce a MapReduce v Hadoop

Dále implementujeme Reducer \equiv obdoba funkce `reduce` v TNL.

```
1  import org.apache.hadoop.io.IntWritable;
2  import org.apache.hadoop.io.Text;
3  import org.apache.hadoop.mapreduce.Reducer;
4
5  import java.io.IOException;
6
7  public class WordCountReducer
8      extends Reducer<Text, IntWritable, Text, IntWritable> {
9
10     private IntWritable result = new IntWritable();
11
12     public void reduce(Text key,
13                       Iterable<IntWritable> values,
14                       Context context)
15         throws IOException, InterruptedException {
16         int sum = 0;
17         for (IntWritable val : values) {
18             sum += val.get();
19         }
20         result.set(sum);
21         context.write(key, result);
22     }
23 }
```

- ▶ Reducer dostává pro jednotlivé klíče množinu hodnot tak, je generoval Mapper. S hodnotami provede redukci.

Paralelní redukce a MapReduce v Hadoop

Nakonec implementujeme `Driver`, který inicializuje zpracování dat.

```
1 import org.apache.hadoop.conf.Configuration;
2 import org.apache.hadoop.fs.Path;
3 import org.apache.hadoop.io.IntWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Job;
6 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
7 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
8
9 public class WordCount {
10     public static void main(String[] args) throws Exception {
11         if (args.length != 2) {
12             System.err.println("Usage: WordCount <input path> <output path>");
13             System.exit(-1);
14         }
15
16         Configuration conf = new Configuration();
17         Job job = Job.getInstance(conf, "word count");
18         job.setJarByClass(WordCount.class);
19         job.setMapperClass(WordCountMapper.class);
20         job.setCombinerClass(WordCountReducer.class);
21         job.setReducerClass(WordCountReducer.class);
22         job.setOutputKeyClass(Text.class);
23         job.setOutputValueClass(IntWritable.class);
24
25         FileInputFormat.addInputPath(job, new Path(args[0]));
26         FileOutputFormat.setOutputPath(job, new Path(args[1]));
27
28         System.exit(job.waitForCompletion(true) ? 0 : 1);
29     }
30 }
```

Paralelní redukce a MapReduce v Hadoop

Příklad pak lze spustit následujícím příkazem v bashi:

```
1  hadoop jar WordCount.jar WordCount /path/to/input /path/to/output
```

Prefix sum

Definition

Mějme pole prvků a_1, \dots, a_n určitého typu a asociativní operaci \oplus . **Inkluzivní prefix sum** je aplikace asociativní operace \oplus na všechny prvky vstupního pole tak, že výsledkem je pole s_1, \dots, s_n stejného typu, pro kterou platí

$$s_i = \bigoplus_{j=1}^i a_j.$$

Exkluzivní prefix sum je definován vztahy $\sigma_1 = 0$ a

$$\sigma_i = \bigoplus_{j=1}^{i-1} a_j,$$

pro $i > 0$.

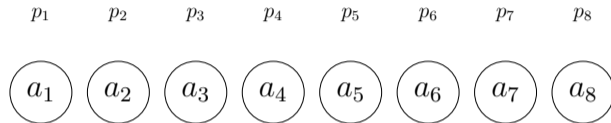
Prefix sum

Poznámka: $\sigma_i = s_i - a_i$ pro $i = 1, \dots, n$.

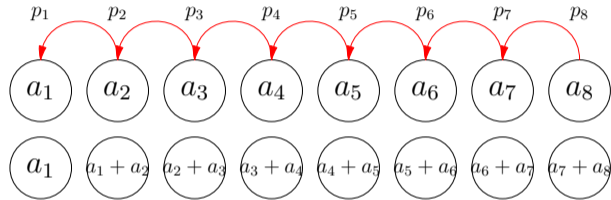
Příklad: Pro jednoduchost budeme uvažovat operaci sčítání. Pak jde o paralelní provedení tohoto kódu:

```
s[ 1 ] = a[ 1 ];  
for( int i = 2; i <= n; i ++ )  
    s[ i ] += a[ i ] + s[ i - 1 ];
```

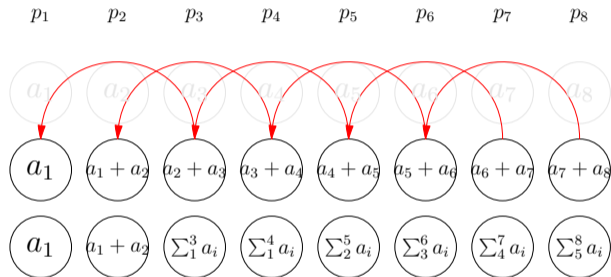
Paralelní prefix sum



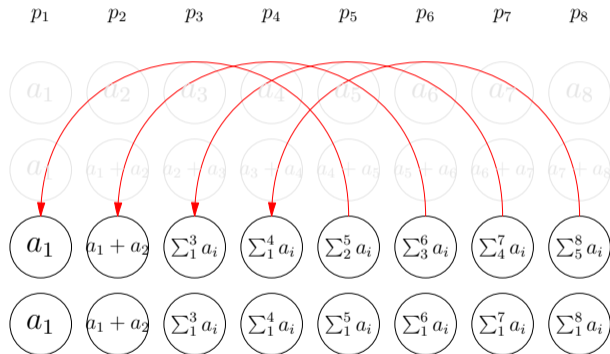
Paralelní prefix sum



Paralelní prefix sum



Paralelní prefix sum



Paralelní prefix sum

- ▶ $T_S(n) = \theta(n)$
- ▶ $T_P(n) = \theta(\log n)$
- ▶ $S(n) = \theta\left(\frac{n}{\log n}\right) = O(n) = O(p)$
- ▶ $E(n) = \theta\left(\frac{1}{\log n}\right)$

Paralelní prefix sum

- ▶ vidíme, že prefix sum má stejnou složitost jako redukce
- ▶ abysme získali nákladově optimální algoritmus, potřebovali bysme na jeden procesor mapovat celý blok čísel
- ▶ nyní je to o trochu složitější

Nákladově optimální paralelní prefix sum

- ▶ volíme $p < n$, pro jednoduchost tak, aby p dělilo n
- ▶ necht' každý procesor má blok A_k pro $k = 1, \dots, p$
- ▶ v něm se sekvenčně napočítá částečný prefix sum

$$s_i^* = \sum_{j \in A_i \wedge j \leq i} a_j.$$

- ▶ definujeme posloupnost S_k pro $k = 1, \dots, p$ tak, že $S_k = s_{i_k}$, kde i_k je index posledního prvku v bloku A_k .
- ▶ napočítáme exkluzivní prefix sum z posloupnosti $S_k \rightarrow \Sigma_k$
- ▶ každý procesor nakonec přičte Σ_k ke svým s_i
- ▶ potom je $T_P(n, p) = \theta \left(\frac{n}{p} + \log p \right)$ stejně jako u redukce

Segmentovaný prefix-sum

- ▶ jde o napočítání několik prefix-sum najednou
- ▶ například

[1, 2, 3], [4, 5, 6, 7, 8]

dá výsledek

[1, 3, 6], [4, 9, 15, 22, 30]

- ▶ úlohu lze zakódovat jako jednu dlouhou řadu s restartovacími značkami
- ▶ sekvenční implementace pak může vypadat takto

Segmentovaný prefix-sum

```
1 void sequentialSegmentedPrefixSum( const int* a,  
2                                   const int* segments,  
3                                   const int n,  
4                                   int* segmentedPrefixSum )  
5 {  
6     segmentedPrefixSum[ 0 ] = a[ 0 ];  
7     for( int i = 1; i < n; i ++ )  
8         if( segments[ i ] == 0 )  
9             segmentedPrefixSum[ i ] =  
10            segmentedPrefixSum[ i - 1 ] + a[ i ];  
11     else  
12         segmentedPrefixSum[ i ] = a[ i ];  
13 }
```

Segmentovaný prefix-sum

- ▶ teoreticky lze segmentovaný prefix-sum implementovat stejně jako normální jen s jinak definovanou operací "sčítání"

$$(s_i, f_i) \oplus (s_j, f_j) = (f_j == 0 ? s_i + s_j : s_j, f_i | f_j),$$

kde $i < j$ a f_i je posloupnost jedniček na začátcích segmentů, jinak samé nuly