

# Linguagem Chapel

Walter Perez Urcia

Universidade de São Paulo  
Instituto de Matemática e Estatística  
Departamento de Ciências da Computação

17 de junho de 2015

# Timeline

## 1 Introdução

## 2 Chapel

- Instalação e configuração
- Tipos de dados
- Semelhanças com outras linguagens
- Data Parallelism
- Task Parallelism
- Multi-locale Parallelism
- The Model and the Problem
- The Integrated Approach

## 3 Bad News: Hardness Results

- Hardness of PP-Partitioning of Haplotype Matrices
- Hardness of PP-Partitioning of Genotype Matrices

## 4 Good News: Tractability Results

- Perfect Path Phylogenies
- Tractability of PPP-Partitioning of Genotype Matrices

## Soma de matrizes

Quantas linhas são necessários para a sua implementação?

# Soma de matrizes

## OpenMP

```
1 #pragma omp shared( A , B , C , chunk ){  
2   #pragma omp for schedule( static , chunk )  
3   for( i = 0 ; i < N ; i++)  
4     for( j = 0 ; j < N ; j++)  
5       C[ i ][ j ] = A[ i ][ j ] + B[ i ][ j ] ;  
6 }
```

# Soma de matrizes

## CUDA

```
1 const int BlockSizeX = 32 ;
2 const int Factor = 4 ;
3 const int BlockSizeY = BlockSizeX / Factor ;
4 __global__ void add( int* C , int* A , int* B , const int N ){
5     int col = blockIdx.x * BlockSizeX + threadIdx.x ;
6     int row = blockIdx.y * BlockSizeX + threadIdx.y * Factor ;
7     for( i = 0 ; i < Factor ; i++)
8         C[ ( row + i ) * N + col ] = A[ ( row + i ) * N + col ] + B[
9             ( row + i ) * N + col ] ;
}
```

## Soma de matrizes

### OpenMPI

```
1 MPI_Comm_size( MPI_COMM_WORLD , &npes ) ;
2 MPI_Comm_rank( MPI_COMM_WORLD , &myrank ) ;
3 if( myrank == ROOT )
4     for( target = 0 ; target < npes ; i++){
5         MPI_Send( A+N*target , N , MPI_INT , target , target ,
6                 MPI_COMM_WORLD ) ;
7         MPI_Send( B+N*target , N , MPI_INT , target , target ,
8                 MPI_COMM_WORLD ) ;
9     }
10 MPI_Recv( myA , N , MPI_INT , ROOT , myrank , MPI_COMM_WORLD , &
11          status ) ;
12 MPI_Recv( myB , N , MPI_INT , ROOT , myrank , MPI_COMM_WORLD , &
13          status ) ;
14 for( i = 0 ; i < N ; i++) myC[ i ] = myA[ i ] + myB[ i ] ;
15 MPI_Send( myC , N , MPI_INT , ROOT , 0 , MPI_COMM_WORLD ) ;
16 if( myrank == ROOT )
17     for( sender = 0 ; sender < npes ; sender++){
18         MPI_Recv( C+N*sender , N , MPI_INT , sender , 0 ,
19                 MPI_COMM_WORLD , &status ) ;
```

## Soma de matrizes

### Linhas

- OpenMP: 6
- CUDA: 9
- OpenMPI: 14

Então, quantas linhas são necessários com Chapel?

## Our formalization of haplotyping.

### Inputs

- A genotype matrix  $G$ .
- The rows of the matrix are individuals / taxa.
- The columns of the matrix are SNP sites / characters.
- The problem is directed: one haplotype is known.
- The input is biallelic: there are only two homozygous states (0 and 1) and one heterozygous state (2).

### Outputs

- A haplotype matrix  $H$ .
- Pairs of rows in  $H$  explain the rows of  $G$ .
- The haplotypes in  $H$  form a perfect phylogeny.



## We can do perfect phylogeny haplotyping efficiently, but ...

### ① Data may be missing.

- This makes the problem NP-complete ...
- ... even for very restricted cases.

#### Solutions:

- Additional assumption like the rich data hypothesis.

### ② No perfect phylogeny is possible.

- This can be caused by chromosomal crossing-over effects.
- This can be caused by incorrect data.
- This can be caused by multiple mutations at the same sites.

#### Solutions:

- Look for phylogenetic networks.
- Correct data.
- Find blocks where a perfect phylogeny is possible.

## How blocks help in perfect phylogeny haplotyping.

- 1 Partition the site set into overlapping contiguous blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 Use dynamic programming for finding the partition.

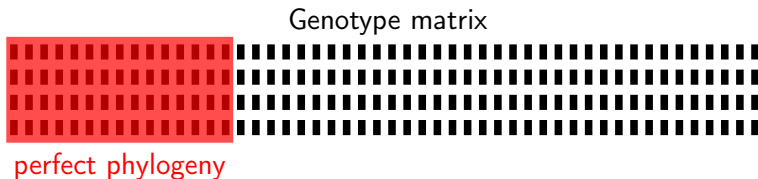
Genotype matrix



no perfect phylogeny

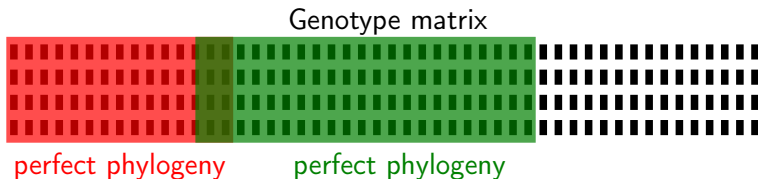
## How blocks help in perfect phylogeny haplotyping.

- 1 Partition the site set into overlapping contiguous blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 Use dynamic programming for finding the partition.



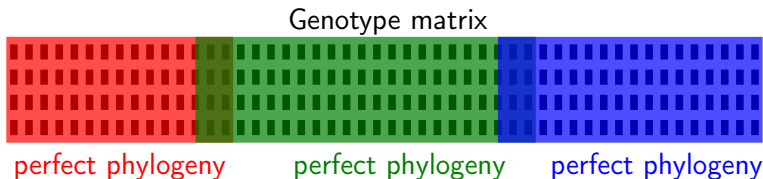
## How blocks help in perfect phylogeny haplotyping.

- 1 Partition the site set into overlapping contiguous blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 Use dynamic programming for finding the partition.



## How blocks help in perfect phylogeny haplotyping.

- 1 Partition the site set into overlapping contiguous blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 Use dynamic programming for finding the partition.



## Objective of the integrated approach.

- 1 Partition the site set into **noncontiguous** blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 **Compute partition while computing perfect phylogenies.**

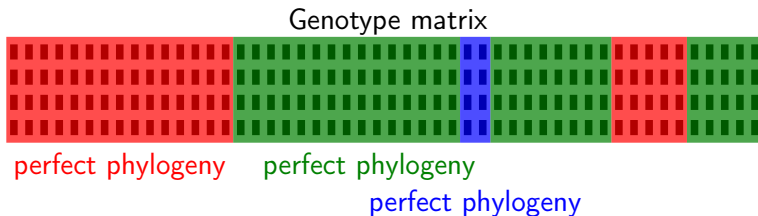
Genotype matrix



no perfect phylogeny

## Objective of the integrated approach.

- 1 Partition the site set into **noncontiguous** blocks.
- 2 Compute a perfect phylogeny for each block and combine them.
- 3 Compute partition while computing perfect phylogenies.



## The formal computational problem.

We are interested in the computational complexity of the function  $\chi_{PP}$ :

- It gets genotype matrices as input.
- It maps them to a number  $k$ .
- This number is minimal such that the sites can be covered by  $k$  sets, each admitting a perfect phylogeny.  
(We call this a **pp-partition**.)



## Finding pp-partitions of haplotype matrices.

We start with a special case:

- The inputs  $M$  are already haplotype matrices.
- The inputs  $M$  do not allow a perfect phylogeny.
- What is  $\chi_{PP}(M)$ ?

### Example

$M$ :

0	0	0	1
0	1	0	0
1	0	0	0
0	1	0	0
1	0	0	0
0	1	0	1
1	1	0	0
0	0	1	0
1	0	1	0

No perfect phylogeny is possible.

## Finding pp-partitions of haplotype matrices.

We start with a special case:

- The inputs  $M$  are already haplotype matrices.
- The inputs  $M$  do not allow a perfect phylogeny.
- What is  $\chi_{PP}(M)$ ?

### Example

$M$ :

0	0	0	1
0	1	0	0
1	0	0	0
0	1	0	0
1	0	0	0
0	1	0	1
1	1	0	0
0	0	1	0
1	0	1	0

Perfect phylogeny

Perfect phylogeny

$$\chi_{PP}(M) = 2.$$

## Bad news about pp-partitions of haplotype matrices.

### Theorem

Finding *optimal pp-partition of haplotype matrices* is equivalent to finding *optimal graph colorings*.

### Proof sketch for first direction.

- 1 Let  $G$  be a graph.
- 2 Build a matrix with a column for each vertex of  $G$ .
- 3 For each edge of  $G$  add four rows inducing the submatrix  $\begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix}$ .
- 4 The submatrix enforces that the columns lie in different perfect phylogenies.



## Implications for pp-partitions of haplotype matrices.

### Corollary

*If  $\chi_{PP}(M) = 2$  for a haplotype matrix  $M$ , we can find an optimal pp-partition in polynomial time.*

### Corollary

*Computing  $\chi_{PP}$  for haplotype matrices is*

- *NP-hard,*
- *not fixed-parameter tractable, unless  $P = NP$ ,*
- *very hard to approximate.*

## Finding pp-partitions of genotype matrices.

Now comes the general case:

- The inputs  $M$  are **genotype matrices**.
- The inputs  $M$  **do not allow a perfect phylogeny**.
- What is  $\chi_{PP}(M)$ ?

### Example

$M$ :

2	2	2	2
1	0	0	0
0	0	0	1
0	0	1	0
0	2	2	0
1	1	0	0

No perfect phylogeny is possible.

## Finding pp-partitions of genotype matrices.

Now comes the general case:

- The inputs  $M$  are **genotype matrices**.
- The inputs  $M$  **do not allow a perfect phylogeny**.
- What is  $\chi_{PP}(M)$ ?

### Example

$M$ :

2	2	2	2
1	0	0	0
0	0	0	1
0	0	1	0
0	2	2	0
1	1	0	0

Perfect phylogeny

Perfect phylogeny

$$\chi_{PP}(M) = 2.$$

## Bad news about pp-partitions of haplotype matrices.

### Theorem

Finding *optimal pp-partition of genotype matrices* is at least as hard as finding *optimal colorings of 3-uniform hypergraphs*.

### Proof sketch.

- 1 Let  $G$  be a 3-uniform hypergraph.
- 2 Build a matrix with a column for each vertex of  $G$ .
- 3 For each hyperedge of  $G$  add four rows inducing the submatrix  $\begin{pmatrix} 2 & 2 & 2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ .
- 4 The submatrix enforces that the three columns do not all lie in the same perfect phylogeny. □

## Implications for pp-partitions of genotype matrices.

### Corollary

*Even if we know  $\chi_{PP}(M) = 2$  for a genotype matrix  $M$ , finding a pp-partition of any fixed size is still*

- NP-hard,
- not fixed-parameter tractable, unless  $P = NP$ ,
- very hard to approximate.



## Automatic optimal pp-partitioning is hopeless, but...

- The hardness results are **worst-case** results for **highly artificial inputs**.
- **Real biological data** might have special properties that make the problem **tractable**.
- One such property is that perfect phylogenies are often perfect **path** phylogenies:  
In HapMap data, in 70% of the blocks where a perfect phylogeny is possible a perfect path phylogeny is also possible.

## Example of a perfect path phylogeny.

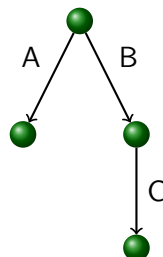
### Genotype matrix

	A	B	C
	2	2	2
$G:$	0	2	0
	2	0	0
	0	2	2

### Haplotype matrix

	A	B	C
	1	0	0
	0	1	1
	0	0	0
$H:$	0	1	0
	0	0	0
	1	0	0
	0	0	0
	0	1	1

### Perfect path phylogeny



## The modified formal computational problem.

We are interested in the computational complexity of the function  $\chi_{\text{PPP}}$ :

- It gets genotype matrices as input.
- It maps them to a number  $k$ .
- This number is minimal such that the sites can be covered by  $k$  sets, each admitting a perfect **path** phylogeny.  
(We call this a ppp-partition.)

## Good news about ppp-partitions of genotype matrices.

### Theorem

*Optimal ppp-partitions of genotype matrices can be computed in polynomial time.*

### Algorithm

- 1 Build the following partial order:
  - Can one column be above the other in a phylogeny?
  - Can the columns be the two children of the root of a perfect path phylogeny?
- 2 Cover the partial order with as few compatible chain pairs as possible. For this, a maximal matching in a special graph needs to be computed.

► The algorithm in action

## Summary

- Finding optimal pp-partitions is **intractable**.
- It is even intractable to find a pp-partition when **just two noncontiguous blocks are known to suffice**.
- For perfect **path** phylogenies, optimal partitions can be computed **in polynomial time**.

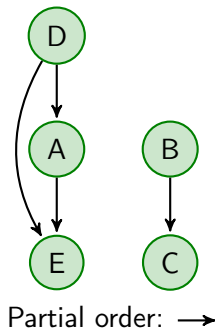
## The algorithm in action.

### Computation of the partial order.

#### Genotype matrix

	A	B	C	D	E
	2	2	2	2	2
G:	0	1	2	1	0
	1	0	0	1	2
	0	2	2	0	0

#### Partial order



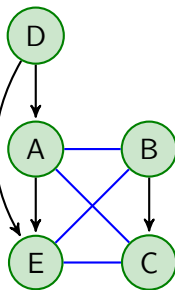
## The algorithm in action.

### Computation of the partial order.

#### Genotype matrix

	A	B	C	D	E
	2	2	2	2	2
G:	0	1	2	1	0
	1	0	0	1	2
	0	2	2	0	0

#### Partial order



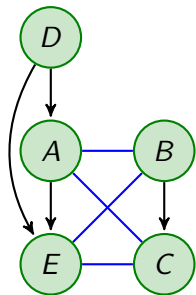
Partial order: →

Compatible as children of root: —

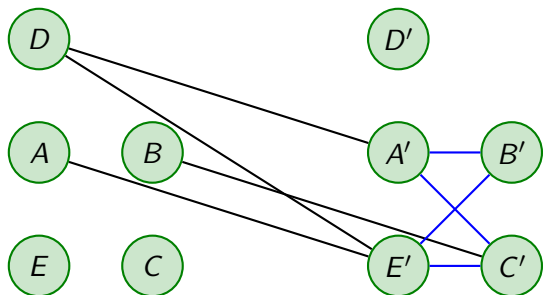
## The algorithm in action.

### The matching in the special graph.

Partial order



Matching graph

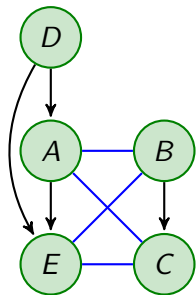




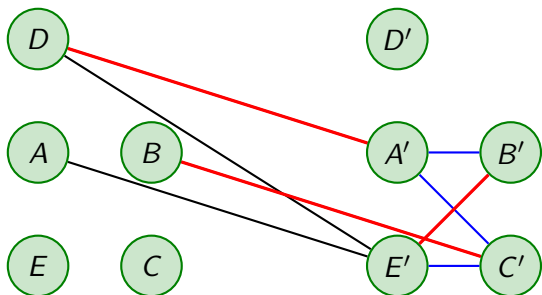
## The algorithm in action.

### The matching in the special graph.

Partial order



Matching graph

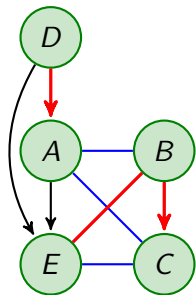


A **maximal matching** in the matching graph

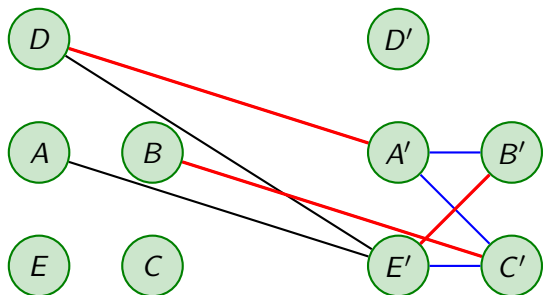
# The algorithm in action.

## The matching in the special graph.

Partial order



Matching graph



A **maximal matching** in the matching graph induces perfect path phylogenies.

## Simple slide with three points shown all at once

- Point 1
- Point 2
- Point 3

## Simple slide with three points shown in succession

- Point 1 (Click “Next Page” to see Point 2)

## Simple slide with three points shown in succession

- Point 1 (Click “Next Page” to see Point 2)
- Point 2

## Simple slide with three points shown in succession

- Point 1 (Click “Next Page” to see Point 2)
- Point 2
- Point 3

## Slide with two columns: items and a graphic

- First item

Insert graphic here

## Slide with two columns: items and a graphic

- First item
- Second item

Insert graphic here



## Slide with two columns: items and a graphic

- First item
- Second item
- ...

Insert graphic here