

Comparing subtelomeric sequences in human genomes in terms of sequence (and methylation) similarity

Jindřich Matuška

Faculty of Informatics, Masaryk University

14. května 2025

Subtelomeres

- Region ca. 500000 bp beside telomeres
- Frequent modifications
- Transposable element

Why do that

- Modifications (duplication) may lead to:
 - creation of new genes
 - diversification of chromosomes

Pipeline

1. Extraction of subtelomeres
2. Similarity by ModDotPlot
3. Further analysis

Extraction of subtelomeres

1. Extraction of buffered sequence from ends (SeqKit, script)
2. Extraction of telomere regions (Seqtk telo)
3. Inversion, cropping subtelomeres to length (bedtools, script, SeqKit)
4. Reversion of end sequences (SeqKit)

Similarity by ModDotPlot

1. Selection of subtelomeres (SeqKit)
2. Split by N-sequences (script, SeqKit)
3. Sort, merge (SeqKit)
4. Similarity analysis (ModDotPlot)
5. Finish?

Further analysis

1. Parsing .bedpe files into heatmap (Python, Matplotlib)
2. Selection of groups (SeqKit)
3. Similarity analysis (ModDotPlot)

Running pipeline

2 slow parts

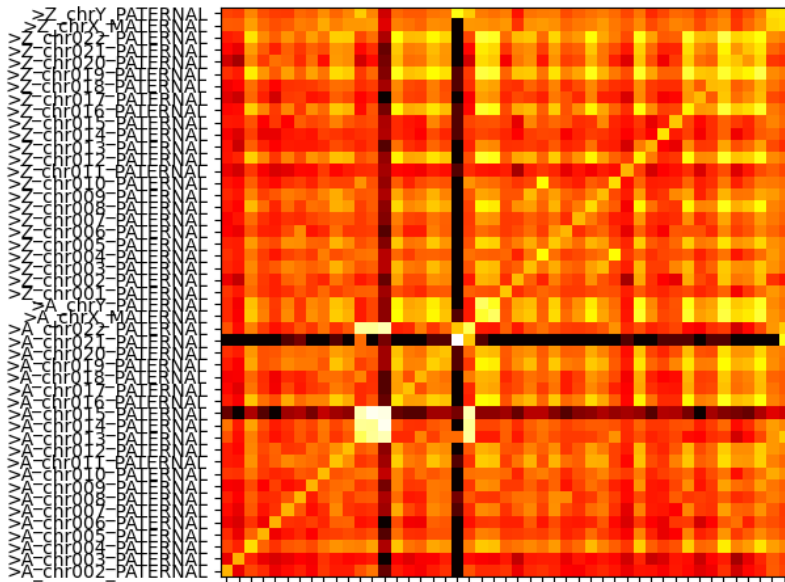
- Extraction of buffered sequences from ends
- ModDotPlot

Otherwise quite fast

Analysis of heatmap

3 groups with high similarity:

- Big tandem repeats
- Singleton chr21-A
- Small similar subsequences

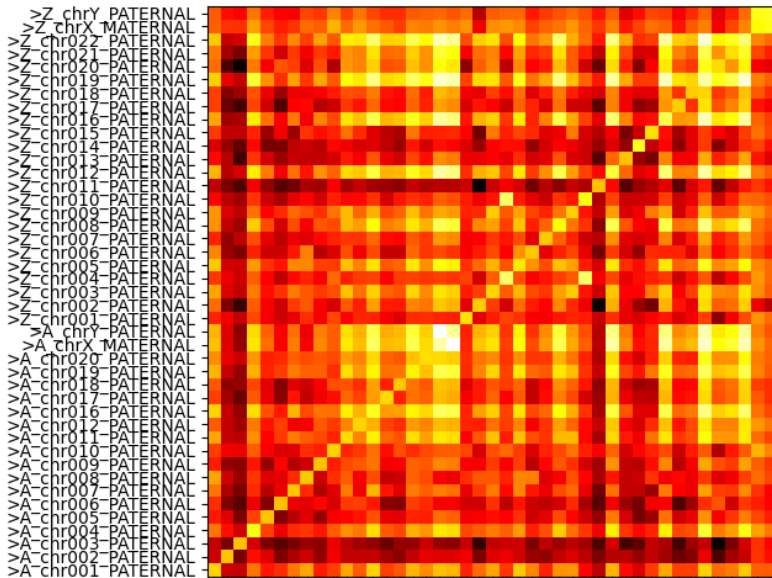


Groups

Big tandem repeats

- chr13-A
- chr14-A
- chr15-A
- chr22-A

Singleton chr21-A



Groups 2

Small similar subsequences

- chr16-A
- chrX-A
- chrY-A
- chr8-Z
- chr12-Z
- chr16-Z
- chr19-Z
- chr20-Z
- chr22-Z

Further work

- Extraction of sequences
- Methylation data
- Multiple sources of data