





Day 1: Computational Pan-Genomics: Status, Promises and Challenges

Jordan Eizenga, Erik Garrison and Tobias Marschall

CPANG18 @ Instituto Gulbenkian de Ciência March, 2018

Rebooting the Human Genome*

"The Human Genome Project was one of mankind's greatest triumphs. But the official gene map that resulted in 2003, known as the "reference genome," is no longer up to the job."

(Antonio Regalado, MIT Technology Review, June 3, 2015)

*https://www.technologyreview.com/s/537916/rebooting-the-human-genome



Future Perspectives in Computational Pan-Genomics

Workshop: 8 – 12 June 2015, Leiden, the Netherlands

Scientific Organizers	 Victor Guryev, ERIBA Groningen Tobias Marschall, Saarland U / MPII Alexander Schönhuth, CWI Amsterdam Fabio Vandin, SDU Odense Kai Ye, St. Louis, Washington U 	
Invited Speakers	 Can Alkan, Bilkent U Ankara Paul De Bakker, UMC Utrecht Valentina Boeva, Institut Curie Paris Francesca Chiaromonte, Penn State U Francesca Ciccarelli, King's College London Evan Eichler, U Washington Eleazar Eskin, UCLA Paul Kersey, EMBL EBI Jan Korbel, EMBL Heidelberg Jens Lagergren, Karolinska Institute Ben Langmead, Johns Hopkins U Veli Mäkinen, U Helsinki Manja Marz, U Jena Paul Medvedev, Penn State U Sven Rahman, U Duisburg-Essen Ben Raphael, Brown U Knut Reinert, FU Berlin Cenk Sahinalp, Simon Fraser U Ole Schulz-Trieglaff, Illumina UK 	

Workshop Paper



Computational pan-genomics: status, promises and challenges

The Computational Pan-Genomics Consortium*

Corresponding author: Tobias Marschall, Center for Bioinformatics as Saarland University and Max Planck Institute for Informatics, Saarland Informatics Campus, Saarbrücken, Germany, Tel.: 496 813 02 70880; E-mail: tmarschall@mpi-inf.mpgde The Computational Plan-Geromics Consortium formed at a workshop held from 8 to 12 june 2015, at the Lorentz Center in Leiden, the Netherlands, with the purpose of providing a cross-disciplinary overview of the emerging discipline of Computational Plan-Geromics. The workshop was organized by Victor Gurvey. Tobias Marschall, Alexander Schohnith (chair, Tabio Vandin, and Val Consortium members are listed at the end of this article

Abstract

Many disciplines, from human genetics and oncology to plant breeding, microbiology and virology, commonly face the challenge of analyzing rapidly increasing numbers of genomes. In case of Homo sapiers, the number of sequenced genomes will approach hundreds of thousands in the next few years. Simply scaling up established bioinformatics pipelines will not be sufficient for leveraging the full potential of such rich genomic data sets. Instead, novel, qualitatively different computational methods and paradigms are needed. We will witness the rapid extension of computational pan-genomics, a new sub-area of research in computational biology. In this article, we generalize existing def nitions and understand a pangenome as any collection of genomic sequences to be analyzed jointly or to be used as a reference. We examine already available approaches to construct and use pan-genomes, discuss the potential beneft so future technologies and methodologies and review open challenges from the vantage point of the above mentioned biological disciplines. As a prominent example for a computational paragement, be particularly highligh the transition from the representation of reference genomes as strings to representations as graphs. We outline how this and other challenges from different application do-

Table of Content

Introduction

- Definition of Computational Pan-Genomics
- Goals of Computational Pan-Genomics
- Applications
 - Microbes, Metagenomics, Viruses, Plants, Human Genetic Diseases, Cancer, Phylogenomics
- Impact of Sequencing Technology on Pan-Genomics
- Data Structures
 - Design Goals
 - Approaches
- Computational Challenges
 - Read Mapping
 - Variant Calling and Genotyping
 - Haplotype Phasing
 - Visualization
 - Data Uncertainty Propagation

Table of Content

Introduction

- Definition of Computational Pan-Genomics
- Goals of Computational Pan-Genomics
- Applications
 - Microbes, Metagenomics, Viruses, Plants, Human Genetic Diseases, Cancer, Phylogenomics
- Impact of Sequencing Technology on Pan-Genomics
- Data Structures
 - Design Goals
 - Approaches
- Computational Challenges
 - Read Mapping
 - Variant Calling and Genotyping
 - Haplotype Phasing
 - Visualization
 - Data Uncertainty Propagation

Term pan-genome popularized in microbiology in 2005.

Definition (gene-based pan-genome)

Term pan-genome popularized in microbiology in 2005.

Definition (gene-based pan-genome)



Term pan-genome popularized in microbiology in 2005.

Definition (gene-based pan-genome)



Term pan-genome popularized in microbiology in 2005.

Definition (gene-based pan-genome)



Pan-Genome Definition

"... and use the term **pan-genome** to refer to any collection of genomic sequences to be analyzed jointly or to be used as a reference. These sequences can be linked in a graph-like structure, or simply constitute sets of (aligned or unaligned) sequences. Questions about efficient data structures, algorithms and statistical methods to perform bioinformatic analyses of pan-genomes give rise to the discipline of computational pan-genomics."

Pan-Genome Definition

"... and use the term **pan-genome** to refer to any collection of genomic sequences to be analyzed jointly or to be used as a reference. These sequences can be linked in a graph-like structure, or simply constitute sets of (aligned or unaligned) sequences. Questions about efficient data structures, algorithms and statistical methods to perform bioinformatic analyses of pan-genomes give rise to the discipline of computational pan-genomics."

Notes

- Not restricted to taxonomic units
- Not restricted to full genomes
- Not tied to graphs
- Intentionally intersects with metagenomics, comparative genomics, and population genetics

Table of Content

Introduction

- Definition of Computational Pan-Genomics
- Goals of Computational Pan-Genomics
- Applications
 - Microbes, Metagenomics, Viruses, Plants, Human Genetic Diseases, Cancer, Phylogenomics
- Impact of Sequencing Technology on Pan-Genomics
- Data Structures
 - Design Goals
 - Approaches
- Computational Challenges
 - Read Mapping
 - Variant Calling and Genotyping
 - Haplotype Phasing
 - Visualization
 - Data Uncertainty Propagation

(High-Level) Goals of Computational Pan-Genomics

- completeness: containing all functional elements and enough of the sequence space to serve as a reference for the analysis of additional individuals,
- stability: having uniquely identifiable features that can be studied by different researchers and at different points in time,
- comprehensibility: facilitating understanding of the complexities of genome structures across many individuals or species,
- efficiency organizing data in such a way as to accelerate downstream analysis.

Table of Content

Introduction

- Definition of Computational Pan-Genomics
- Goals of Computational Pan-Genomics
- Applications
 - Microbes, Metagenomics, Viruses, Plants, Human Genetic Diseases, Cancer, Phylogenomics
- Impact of Sequencing Technology on Pan-Genomics
- Data Structures
 - Design Goals
 - Approaches
- Computational Challenges
 - Read Mapping
 - Variant Calling and Genotyping
 - Haplotype Phasing
 - Visualization
 - Data Uncertainty Propagation

Approaches I

Haplotype 1 CAAATAAGGCTTGGAAATTTACCCGCTCCTGCCCGCGTCTGGAGTTCACCCGCTCCTGCCCCGCGTATTATATTCCAACTCTCG

- Haplotype 2 CAAATAAGCCTTGGAAATTTACCCGCTCCTGCCCGCGTCTGGAGGTTCTATTATATTCCAACTCTCTG

(a) Unaligned sequences



Approaches II



12

Approaches III



Course Objectives

- Understand Pan-Genomics concepts and appreciate the limitations of linear reference genomes,
- Learn to build corresponding analysis pipelines in the VG framework,
- We focus on putting you in a position to design (and debug) pipelines tailored to your use case,
- Some of the practicals address serious research questions (that we do not completely solve in class), the aim is rather to give you the tools to attack them.

Getting to know each other

- What is your background?
- What are your expectations?
- Which data / use cases do you work on (or will work on)?

Practicals

- This is a hands-on course. We collect praticals in this git repository: https://github.com/Pfern/PANGenomics
- We add new practicals for each day in the course of this week
- Ppracticals are a starting point for exploring what VG can do. So please don't only copy paste commands, but try to understand their meaning, modify them, etc.
- Help each other
- Present your results in class
- Don't hesitate to send pull requests ;)

Table of Content

Introduction

- Definition of Computational Pan-Genomics
- Goals of Computational Pan-Genomics
- Applications
 - Microbes, Metagenomics, Viruses, Plants, Human Genetic Diseases, Cancer, Phylogenomics
- Impact of Sequencing Technology on Pan-Genomics
- Data Structures
 - Design Goals
 - Approaches
- Computational Challenges
 - Read Mapping
 - Variant Calling and Genotyping
 - Haplotype Phasing
 - Visualization
 - Data Uncertainty Propagation

Applications



Application Domain: Microbes

- Pan-genomics at the gene level: established workflows and mature software are available
- For a number of microorganisms, pan-genome sequence data is already available
- Microbial pan-genomes support comparative genomics studies (especially given horizontal gene exchange)
- Genome-wide association studies (GWAS) for microbes is an emerging field



Application Domain: Metagenomics

- Metagenomics: set of genomic sequences co-occurrence in an environment
- Questions: taxonomic composition of the sample, presence of certain gene products or whole pathways, and determining which genomes these functional genes are associated with.
- Pan-Genome data structures present the chance to reveal common adaptations to the environment as well as co-evolution of interactions.



Application Domain: Viruses

- Reliable viral haplotype reconstruction is not fully solved
- Patient's viral pan-genome → diagnosis, staging, and therapy selection
- Virus-host interactions: pan-genome structure of a viral population to be directly compared with that of a susceptible host population



Application Domain: Plants

- Large-scale genomics projects completed / under way: Arabidopsis thaliana, rice, maize, sorghum, and tomato
- Plant genomes are large, complex (containing many repeats) and often polyploid
- Having a pan-genome available for a given crop that includes its wild relatives provides a single coordinate system to anchor all known variation and phenotype information



Application Domain: Human Genetic Diseases

- Numerous genes have been successfully mapped for rare monogenic diseases
- Common diseases ← GWAS ← imputation ← catalogs of human genome variation, their linkage disequilibrium (LD) properties
- Pan-Genomics can help to achieve this, especially in difficult genomic regions



Application Domain: Cancer

- Improved detection of somatic mutations, through improved quality of read mapping to polymorphic regions
- Somatic pan-genome describing the general somatic variability in the human population, → accurate baseline for assessing the impact of somatic alterations.
- Vision: personal cancer pan-genome to be built for each tumor patient: single-cell data, haplotype information, sequencing data from circulating tumor cells and DNA, etc.



Application Domain: Phylogenomics

- Computational pan-genomics: repidly extrcat evolutionary signals, such as gene content tables, sequence alignments of shared marker genes, genome-wide SNPs, or internal transcribed spacer (ITS) sequences
- Move beyond only using the best aligned, and most well behaved residues of a multiple sequence alignment (often used in "traditional" phylogenomics)



Operations on Pan-Genomes

