# Next Generation Sequencing
# and
# Bioinformatics Analysis Pipelines

Adam Ameur

National Genomics Infrastructure

SciLifeLab Uppsala

adam.ameur@igp.uu.se

# Today's lecture

- Data analysis and management at NGI/SciLifeLab

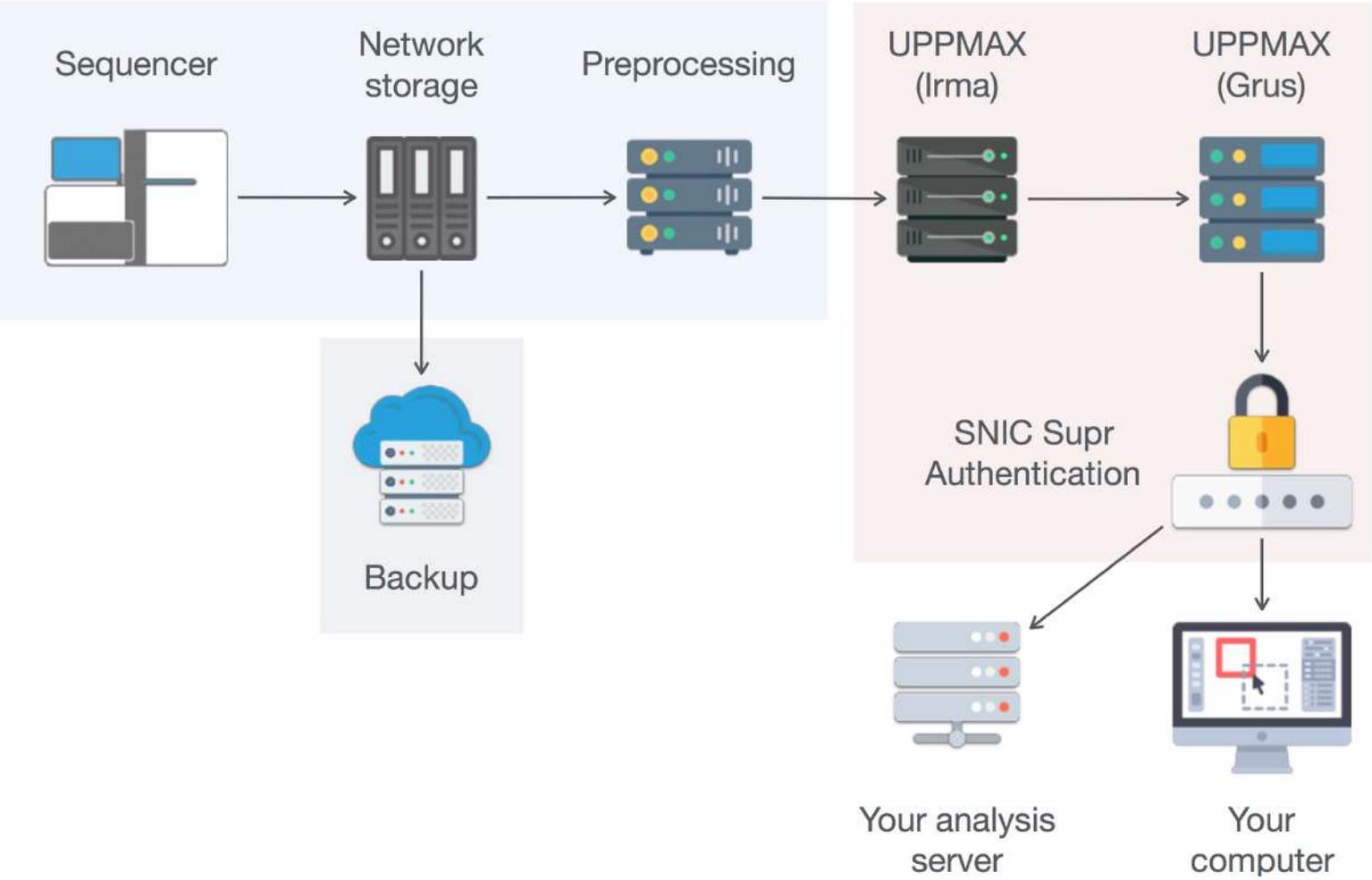- Human whole genome sequencing

- The Earth Biogenome Project

- Other R&D activities at NGI

# NGI Data Handling

# Analysis pipelines

- Initial data analysis for major applications:

  - **Mapping:** Align sequences to a reference genome

  - **SNV calling:** Detect genetic variants

  - **RNA-seq:** Quantify gene expression

  - *De novo* **assembly:** Generate new reference genomes

  - **and more…**

- Analysis requirements: Automated, reliable, easy to run, reproducible

# nf-core

- A community effort to collect a curated set of Nextflow analysis pipelines

- GitHub organisation to collect pipelines in one place

- No institute-specific branding

- Strict set of guideline requirements

https://nf-co.re

## nature biotechnology

Correspondence | Published: 13 February 2020

### The nf-core framework for community-curated bioinformatics pipelines

Philip A. Ewels, Alexander Peltzer, Sven Fillinger, Harshil Patel, Johannes Alneberg, Andreas Wilm, Maxime Ulysse Garcia, Paolo Di Tommaso & Sven Nahnsen ✉

**Phil Ewels, NGI Sthlm**

# Example pipeline - Sarek

GitHub  https://github.com/SciLifeLab/Sarek

Sarek

- Tumour/Normal pair WGS analysis based on GATK best practices

  - SNPs, SNVs and indels

  - Structural variants

  - Heterogeneity, ploidy and CNVs

- Works with regular WGS and Exome data too

| Manta | MuTect1 |
| ASCAT | MuTect2 |
| | Strelka |
| | FreeBayes |
| | GATK HaplotypeCaller |

NATIONAL GENOMICS INFRASTRUCTURE

NBIS NATIONAL BIOINFORMATICS INFRASTRUCTURE SWEDEN

Barntumörbanken

# Quality control

- Every project has some level of quality control checks
  - Technical run performance
  - Read length distribution
  - Sequencing quality


- Analysis pipelines give application-specific QC


- Reporting done using MultiQC (Illumina projects)

# Multi QC example

# Data delivery via GRUS

- GRUS is an UPPMAX tool for NGI data delivery

    - NGI creates a SNIC Supr "delivery project" for each NGI sequencing project

    - Project PI and contact person is given access

    - Email sent with project ID and instructions

- Grus is for secure short-term storage only!

    - Requires two-factor authentication

# NGI Research & Development projects

- For some projects, NGI allocates additional resources for development

  - New applications where we see the need to develop a pipeline

  - Construction of reference datasets and resources

  - Strategic collaborative projects

# Example: The SweGen project

- A whole-genome resource for researchers and clinical labs



*From SweGen release party on Oct 19th 2016*

# SweGen: 1000 Swedish Whole Genomes

- What can the SweGen dataset be used for?

    - Look up genetic variant frequencies

    - Use as matched controls

    - Study population genetics

    - Study human evolutionary history

High demand for the data from many different groups:
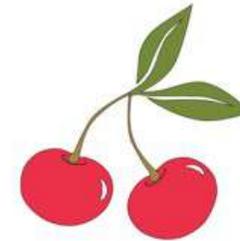
➔ Make the data available as **quickly** and **openly** as possible!

# Deciding on a cohort to use

- The Swedish Twin Registry:

  - Inclusion based on twinning

  - Distribution like population density

  - General population-prevalence of disease

  - 10,000 individuals have been analysed with SNP arrays


- **Identify 1,000 individuals based on genetic structure and diversity across Sweden**

# Selecting 1000 individuals based on PCA

PCA of European samples from the 1000 Genomes Project and 10,000 Swedish samples

# Whole Genome Sequencing

- 30X Illumina WGS generated for all 1,000 individuals

- Sequencing done both at NGI Sthlm and NGI Uppsala

- All 1,000 samples completed in September 2016

# Data analysis pipeline

- NGI pipeline developed for mapping and variant calling



- About 100Gb data generated, and 2 million CPU hours used…

- This pipeline has become standard for all WGS projects at NGI

# Making data available



SweGen Variant Frequency Dataset

This dataset contains whole-genome variant frequencies for 1000 Swedish individuals generated within the SweGen project. The frequency data is intended to be used as a resource for the research community and clinical genetics laboratories.

Please note that the 1000 individuals included in the SweGen project represent a cross-section of the Swedish population and that no disease information has been used for the selection. The frequency data may therefore include genetic variants that are associated with, or causative of, disease.

We request that any use of data from the SweGen project cite this article in the European Journal of Human Genetics.

Individual positions in the genome can be viewed using the Beacon or Graphical Browser. To download the variant frequency file you need to register.

A high confidence set of HLA allele frequencies is available for download under Dataset Access. For a detailed description of the SweGen HLA analysis, please see this bioRxiv preprint.

More information          Beacon          Graphical Browser

- Aggregated frequencies available from: *swefreq.nbis.se*

- Possible to access individual genotype data through Uppmax/Bianca

# SweGen: a resource for collaboration

- Over 90 publications have made use of the SweGen dataset

### Discovery of Novel Sequences in 1,000 Swedish Genomes

Jesper Eisfeldt[*,1,2,3] Gustaf Mårtensson,[4] Adam Ameur[5], Daniel Nilsson[1,2,3] and Anna Lindstrand[1,3]

[1]Department of Molecular Medicine and Surgery, Center for Molecular Medicine, Karolinska Institute, Stockholm, Sweden
[2]Science for Life Laboratory, Karolinska Institutet Science Park, Solna, Sweden

CLINICAL RESEARCH ARTICLE

### Cytokine Autoantibody Screening in the Swedish Addison Registry Identifies Patients With Undiagnosed APS1

Daniel Eriksson,[1,2] Frida Dalin,[1,3] Gabriel Nordling Eriksson,[4] Nils Lan
Matteo Bianchi,[5] Åsa Hallgren,[1,3] Per Dahlqvist,[6] Jeanette Wahlberg,
Olov Ekwall,[10,11] Ola Winqvist,[12] Sergiu-Bogdan Catrina,[4] Johan Rön
Swedish Addison Registry Study Group, Anna-Lena Hulting,[4] Kerstin Lin
Mohammad Alimohammadi,[15] Eystein S. Husebye,[1,16,17,18] Per Morten K
Gerli Rosengren Pielberg,[5] Sophie Bensing,[2,4] and Olle Kämpe[1,2,3,18]

Letter to the Editors-in-Chief

### Prevalence and in silico analysis of missense mutations in the PROS1 gene in the Swedish population: The SweGen dataset

Bengt Zöller

### A rare regulatory variant in the MEF2D gene affects gene regulation and splicing and is associated with a SLE sub-phenotype in Swedish cohorts

Fabiana H. G. Farias, Johanna Dahlqvist, Sergey V. Kozyrev, Dag Leonard, Maria Wilbe, Sergei N. Abramov, Andrei Alexsson, Gerli R. Pielberg, Helene Hansson-Hamlin, Göran Andersson, Karolina Tandre, Anders A. Bengtsson, Christopher Sjöwall, Elisabet Svenungsson, Iva Gunnarsson, Solbritt Rantapää-Dahlqvist, Ann-Christine Syvänen, Johanna K. Sandling, Maija-Leena Eloranta, Lars Rönnblom & Kerstin Lindblad-Toh

- … but also, SweGen is used in clinical routine diagnostics

# Earth Biogenome Project

# EBP – Data management and analysis

- Over the coming years, many new species will be sequenced

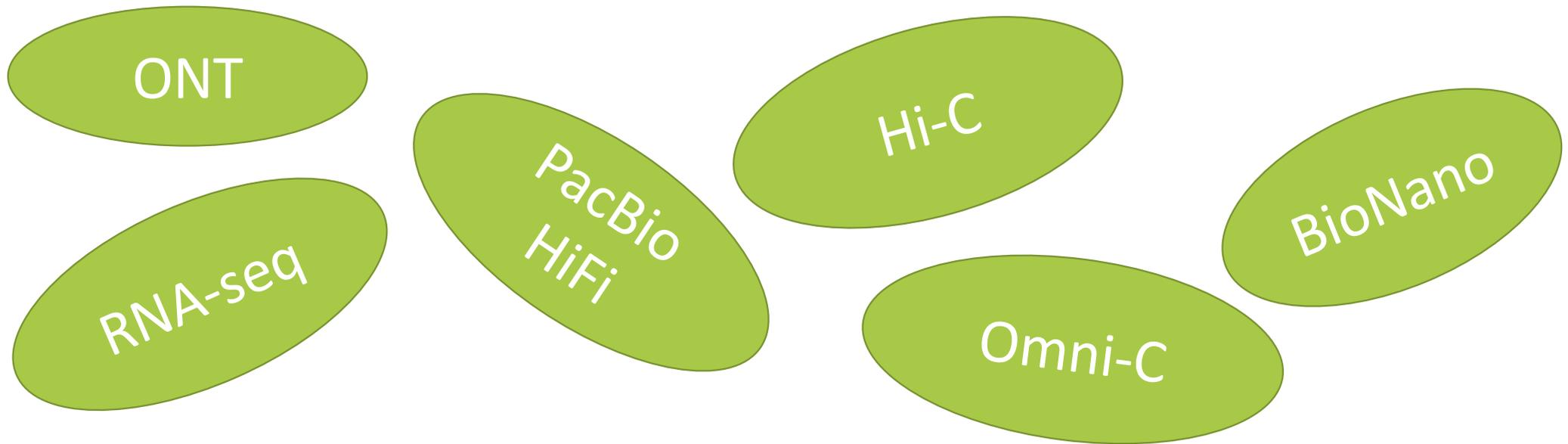- A combination of different instruments and technologies will be used



- We need good strategies for data analysis and management!

# Choice of technology

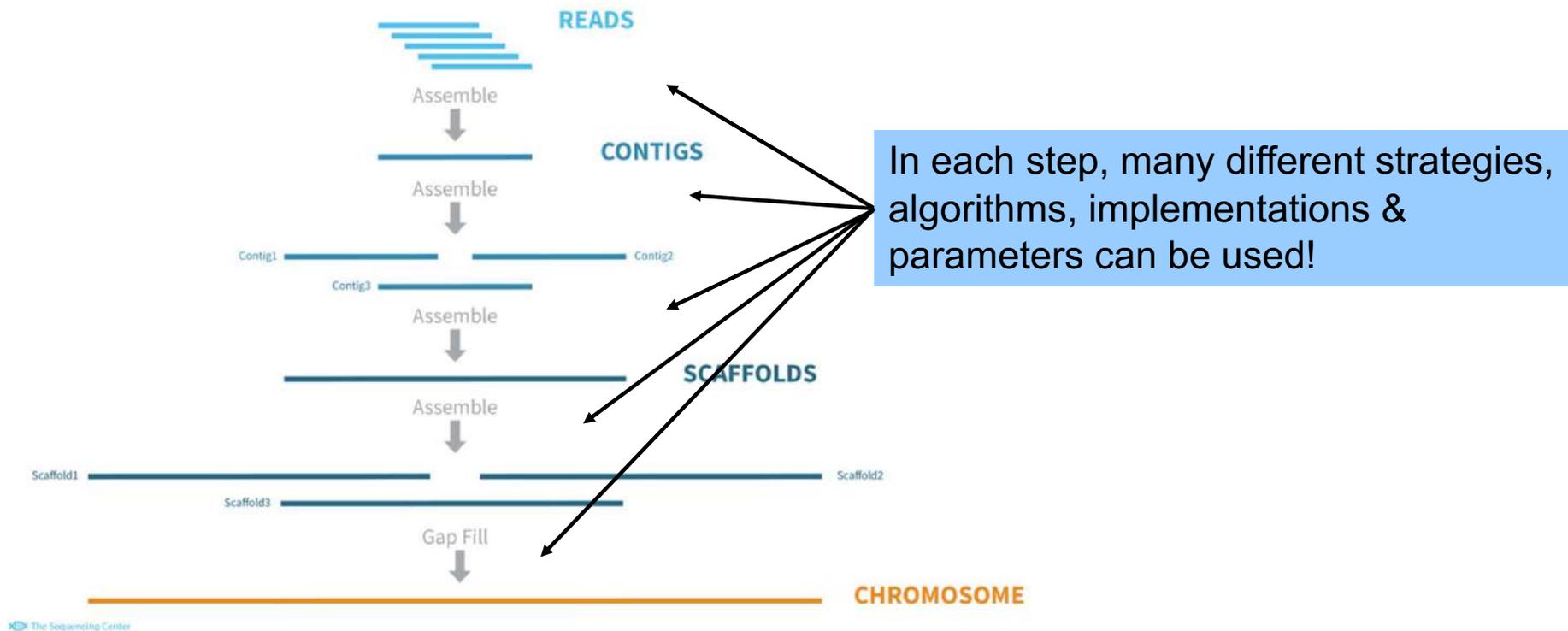- Make sure sequencing is done using the best technology combination

ONT

RNA-seq

PacBio HiFi

Hi-C

Omni-C

BioNano

- This is changing all the time, and lots of different options exist

- The choice will have a big impact on the downstream analysis!

# Genome assembly

- Apply analysis pipelines to generate high-quality genome assemblies



In each step, many different strategies, algorithms, implementations & parameters can be used!

- A challenge for NGI/SciLifeLab is to give best-practice guidelines!

# Genome annotation

- Once the assembly is generated, it needs to be annotated!

- Annotation usually means to find out where genes are located

**Annotation using computational methods**

**Annotation using RNA-sequencing**



- We prefer RNA-sequencing, but still annotation can be challenging!

# Data deposition

- Important to deposit the final assembly in public repositories!



- There is a need to develop an interface to international databases

# EBP – A collaborative project

- A lot of challenges ahead of us to establish EBP analyses in Sweden

- … but the good news is that this is a community effort



- There will likely be a lot or opportunities to collaborate!

# Other R&D activities at NGI

- We have a joint R&D group for all SciLifeLab genomics facilities

- Aim: to test new applications and eventually include as new service

# Examples of technologies evaluated during 2020

**MGI sequencing**

**PacBio Sequel IIe**

**Mission Bio – Tapestri**

# "Milestones and Achievements" 2021

- **New cluster and data delivery system**

- **Olink Explore**

- **Further development of NF-core**

- **Single-cell long read sequencing**

- **Spatial transcriptomics**

If you have other R&D ideas then let us know!

# Olink Explore: measuring proteins by NGS

## What is Olink Explore?

Olink® Explore 1536/384 is a high-multiplex, high-throughput protein biomarker platform that uses Proximity Extension Assay (PEA) technology coupled to an innovative new readout methodology based on Next Generation Sequencing (NGS) using the following Illumina® instruments:

- NovaSeq 6000
- NextSeq 550
- NextSeq 2000



- NGI is the first service provider in Europe for this application!

# CRISPR-Cas9 genome editing

- We are developing new long-read tools to detect "off-target" mutations



*Höijer et al, Genome Biology 2020*

# COVID-19 sequencing

- We are evaluating different methods and sequencing protocols…

**SARS-Cov-2 sequencing (mutations)**

**Human re-sequencing (host response)**



- If you are interested to discuss COVID-19 projects, then let us know!

# Long-read single-cell sequencing

- The first reported WGS of a human single cell with long reads!





*Image from https://en.wikipedia.org/wiki/T_cell*

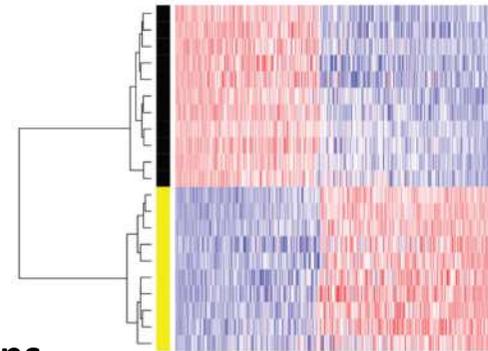- We are also testing new protocols for long-read single-cell RNA sequencing
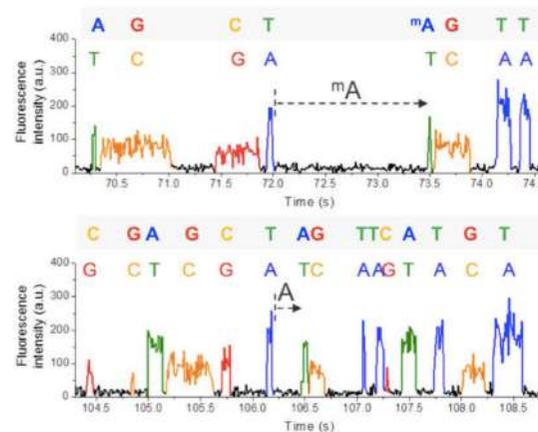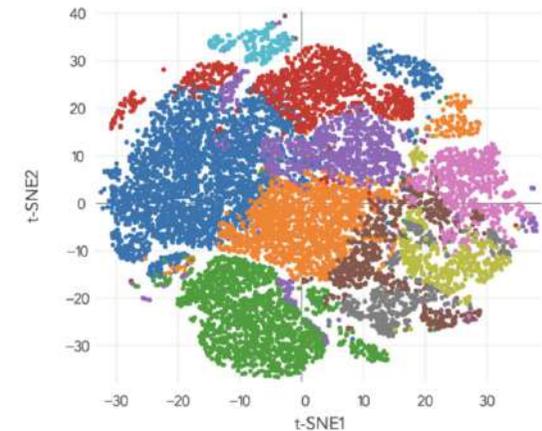
# Many topics I have not covered…

**ChIP-seq analysis**

**RNA-Seq**

**DNA base modifications**

**Single cell RNA-Seq**

- Simply too much to talk about in just one lecture…

# Thanks for your attention!