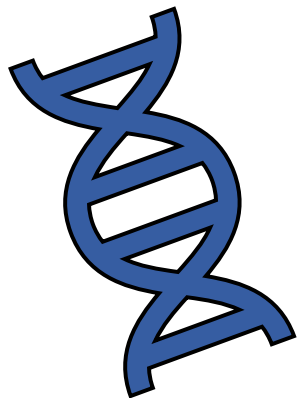
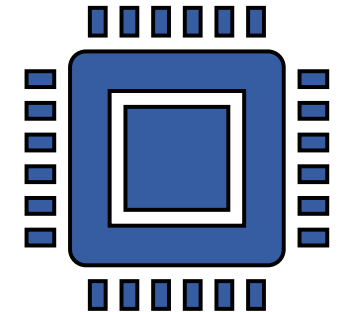


# An Overview of Bioinformatics

BERDC Special Topics Talk 17



## DaCCoTA

DAKOTA COMMUNITY COLLABORATIVE  
ON TRANSLATIONAL ACTIVITY

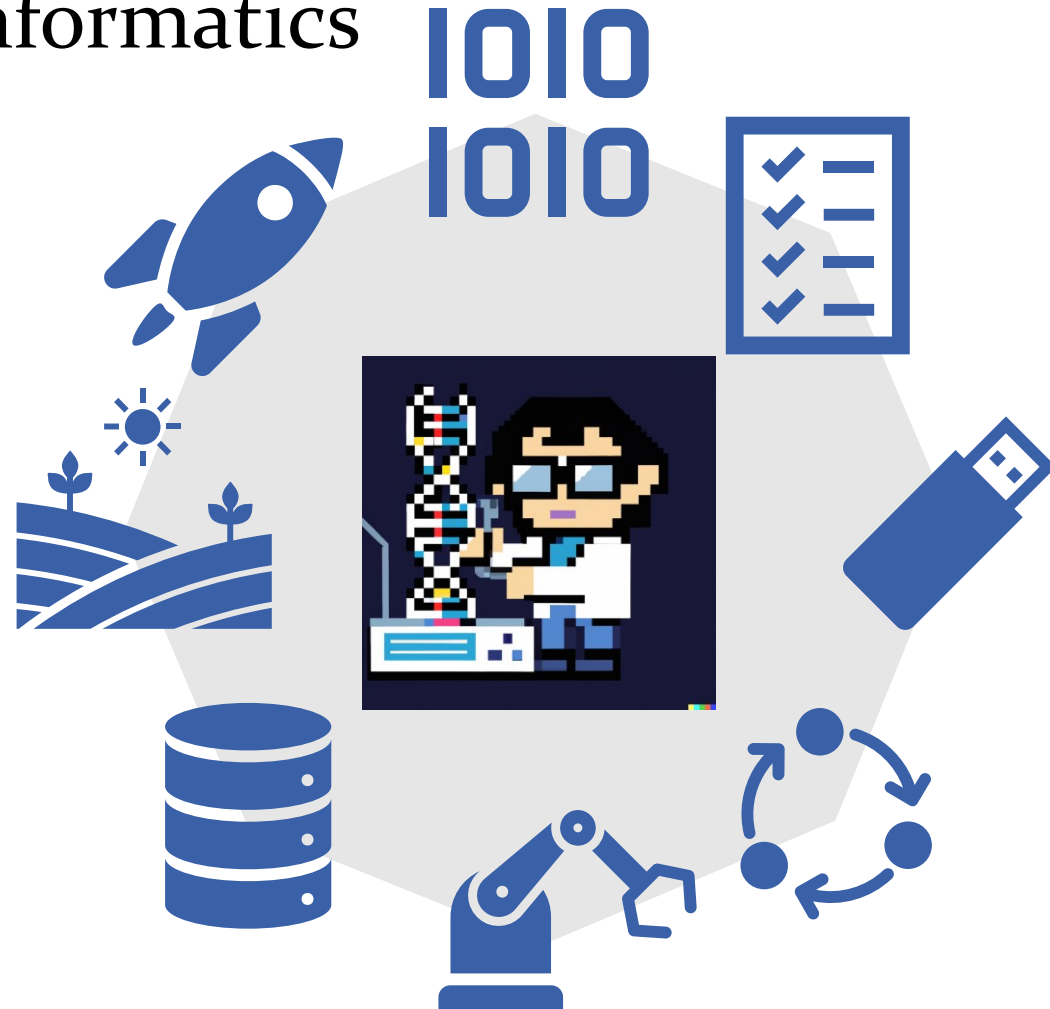
Dr. Mark Williamson

Biostatistics, Epidemiology,  
and Research Design Core

# Introduction

**Goal:** give an overview of Bioinformatics

- 🖨️ Definition & History
- 🖨️ Field & Topics
- 🖨️ Data
- 🖨️ Methods
- 🖨️ Machines & Software
- 🖨️ Databases
- 🖨️ Example Applications
- 🖨️ The Future



# Definition & History

*“the collection, classification, storage, and analysis of biochemical and biological information using computers especially as applied to molecular genetics and genomics” –Merriam-Webster*

- 1960-1970: Protein analysis origins and Margaret Dayhoff [1]
- 1970-1980: Shift from protein to DNA analysis
- 1980-1990: Advances in biology and computing
- 1990-2000: Genomics and the internet
- 2000-2010: High-throughput bioinformatics (2<sup>nd</sup> gen sequencing)
- 2010-Today: Big data, role of bioinformaticians, systems biology

**A** ANALYSIS AND COMPUTATION CENTER

**B** Sequence alignment program output

**C** Cardboard box with handwritten labels: PROT, GDRAN, VER GDRAN, PROT, FAST VER

**D** Protein sequence alignment diagram:

```

  Thr-His-Glu-Cys [Peptide]
  Glu-Cys-Ala-Thr [Peptide]
  Lys-Thr-His [Peptide]
  Met-Ile-Lys [Peptide]
  -----
  Met-Ile-Lys-Thr-His-Glu-Cys-Ala-Thr [Protein]
  
```

**A** match +5 mismatch -4 gap -1

		A	T	C	G
	0	0	0	0	0
A	0	5	-1	-1	-1
T	0	4	10	9	8
G	0	3	9	8	14

**B** Sequence alignment transition diagram:

```

  i-1, j-1  i, j-1
   |         |
   v         v
  i-1, j    i, j
  
```

Score (i,j) = max

**C** Best Alignment : ATCG (Score = 38)

```

  |||
  AT G
  
```

- Score (i-1, j-1) + Match / Mismatch
- Score (i, j-1) + gap
- Score (i-1, j) + gap

# Definition & History

*“the collection, classification, storage, and analysis of biochemical and biological information using computers especially as applied to molecular genetics and genomics” –Merriam-Webster*

1960-1970: Protein analysis origins and Margaret Dayhoff [1]

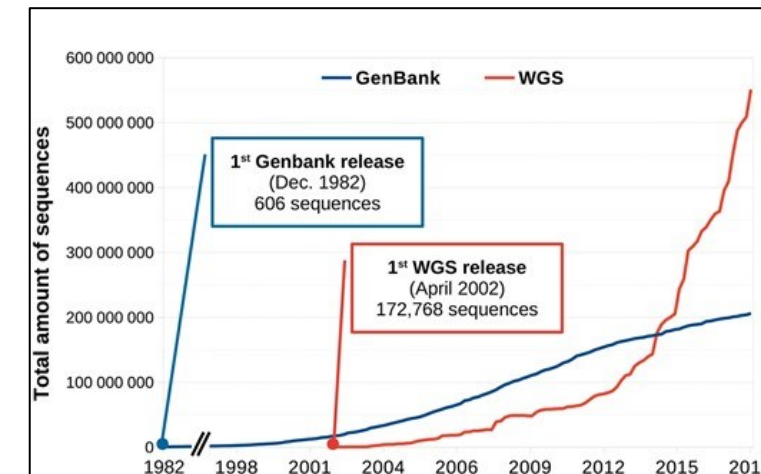
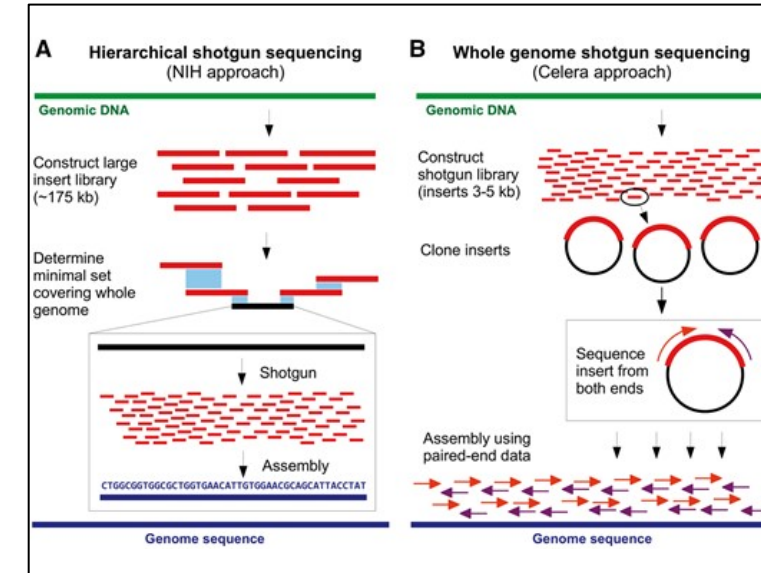
1970-1980: Shift from protein to DNA analysis

1980-1990: Advances in biology and computing

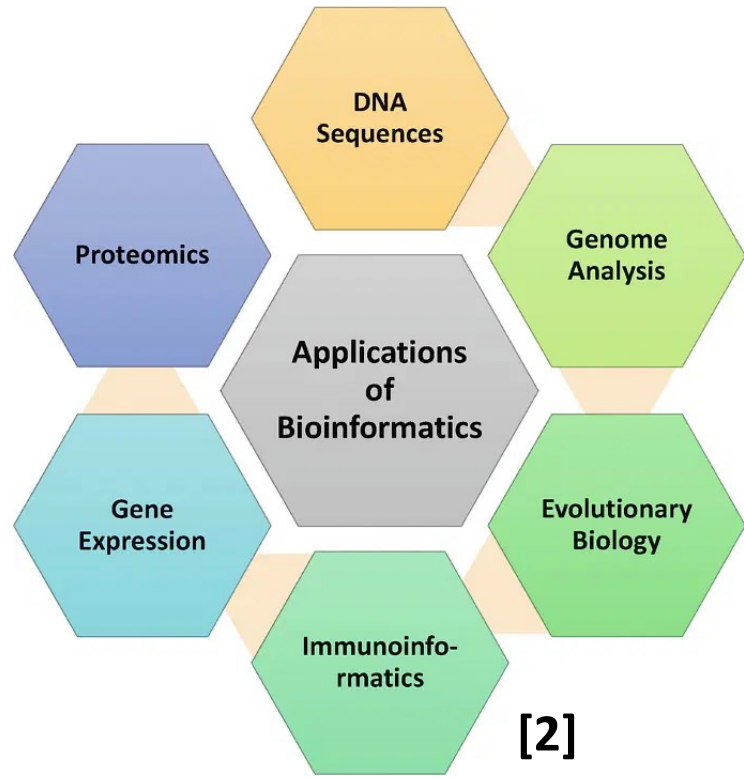
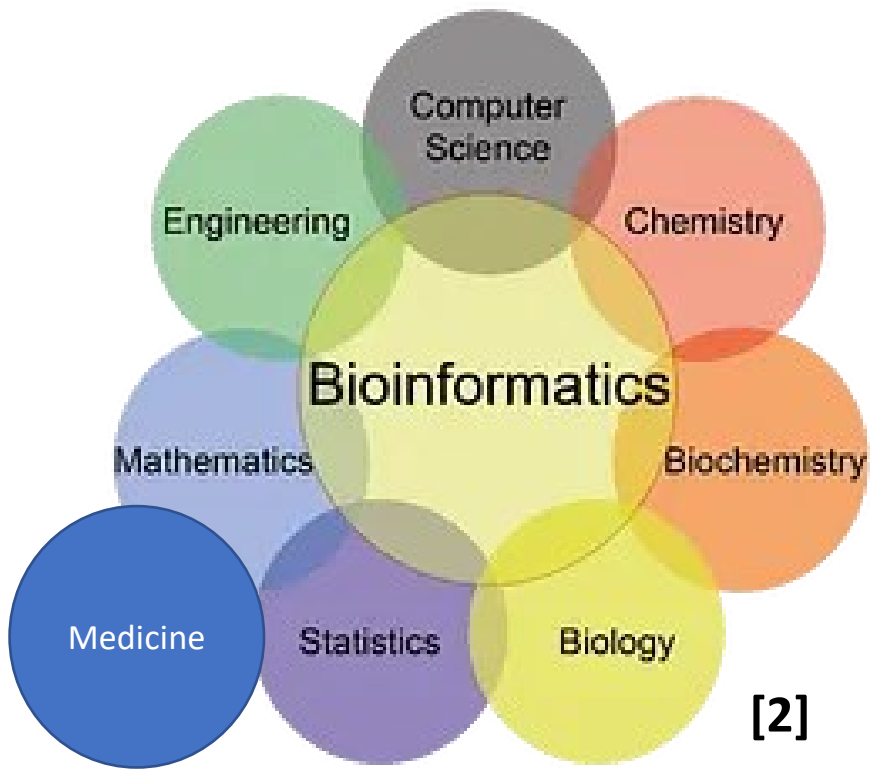
1990-2000: Genomics and the internet

2000-2010: High-throughput bioinformatics (2<sup>nd</sup> gen sequencing)

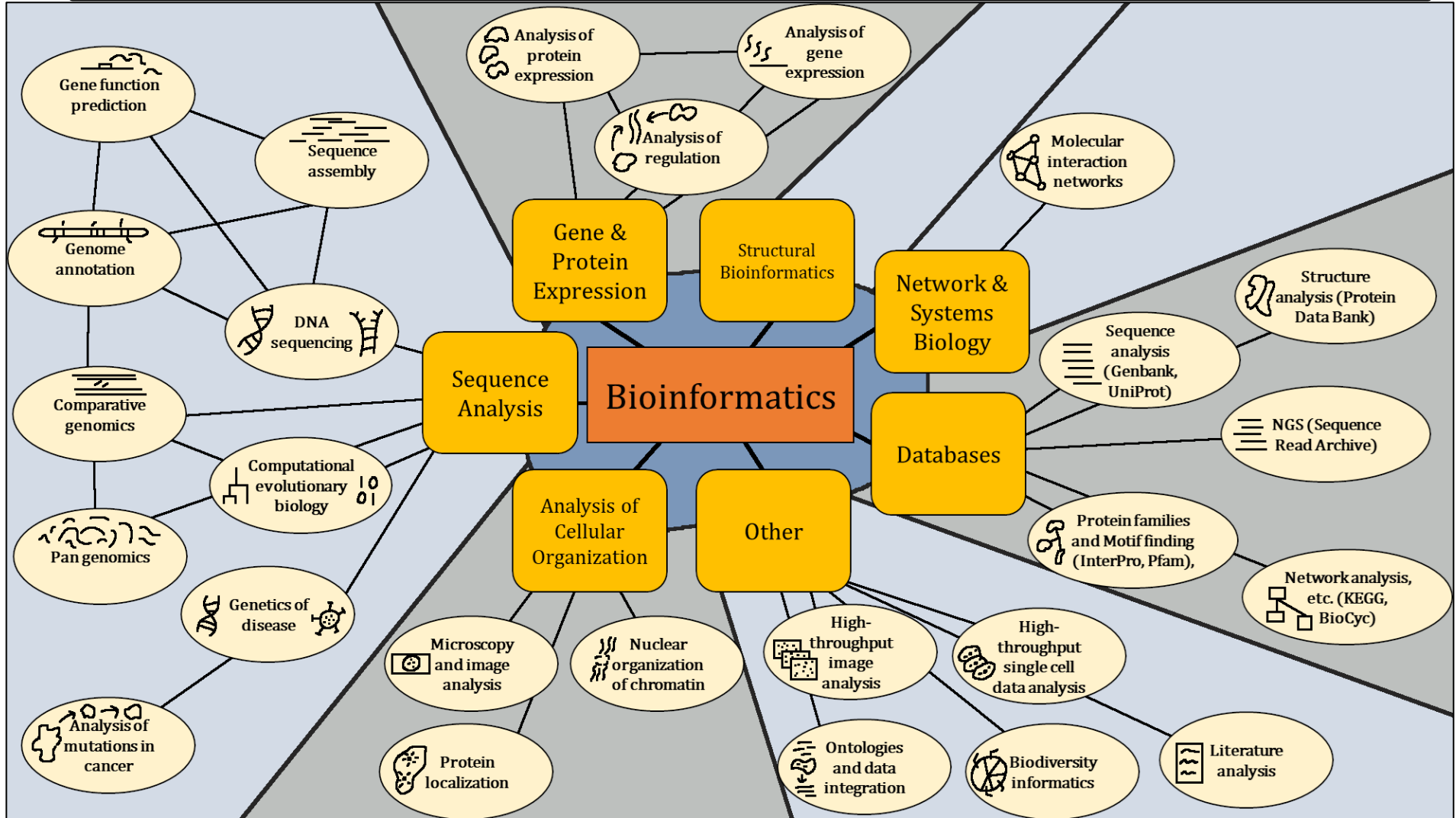
2010-Today: Big data, role of bioinformaticians, systems biology



# Fields & Topics



# Fields & Topics



# Data

## Biological sources

- DNA: sequences, genes, genomes, metagenomes

- RNA and Proteins

- Others [4]

- Metadata

- Images

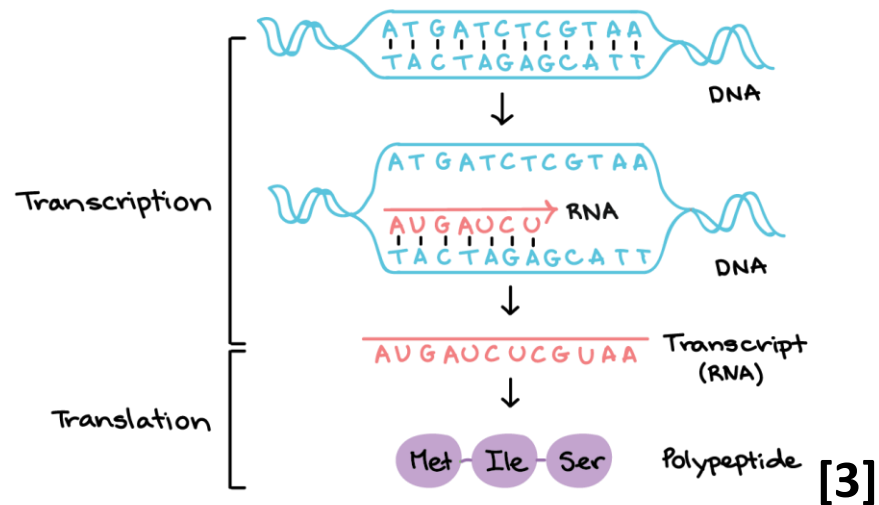
- Etc.

## File formats [5]

- FASTA and FASTQ

- Alignments (BAM, SAM, CRAM)

- Other (VCF, GFF, GTF, BED, Tar.gz, PDB, PED, MAP, CSV, JSON)



>Streptomyces sp. PAMC  
26508 | GCF\_000364805.1 | NC\_021055.1 | Chromosome:  
ANONYMOUS | 3355624..3357149 -  
AAAGGAGGTGATCCAGCCGCACCTTCCGGTACGGCTACCT  
TGTTACGACTTCGTCCCAATCGCCAGTCCCACCTTCGAC  
AGCTCCCTCCACAAGGGGTTGGGCCACCGGCTTCGGGTG  
TTACCGACTTTCGTGACGTGACGGGCGGTGTGTACAAGG  
CCCGGAACGTATTCACCGCAGCAATGCTGATCTGCGATT  
ACTAGCAACTCCGACTTCATGGGGTCGAGTTGCAGACCC  
CAATCCGAAGTACGACCGGCTTTTTGAGATTCGCTCCGCCT  
CGCGGCATCGCAGCTCATTGTACCGGCCATTGTAGCAC  
GTGTGCAGCCCAAGACATAAGGGGCATGATGACTTGACG  
TCGTCCCCACCTTCTCCGAGTTGACCCGGCAGTCTCCT  
GTGAGTCCCCATCACCCGAAGGGCATGCTGGCAACACA  
GAACAAGGGTTGCGCTCGTTGCGGGACTTAACCCAACATC  
TCACGACACGAGCTGACGACAGCCATGCACCACCTGTATA  
CCGACCACAAGGGGGGCACCATCTCTGATGCTTCCGGT

# Methods

🔧 Pipelines

🔧 DNA Sequencing

🔧 Amplicon sequencing: 16/18S, single copy, etc.

🔧 Genome sequencing: reference based, *de novo*

🔧 Annotation (what's there)

🔧 Comparison (what's different between samples/organisms)

🔧 Gene and Protein expression

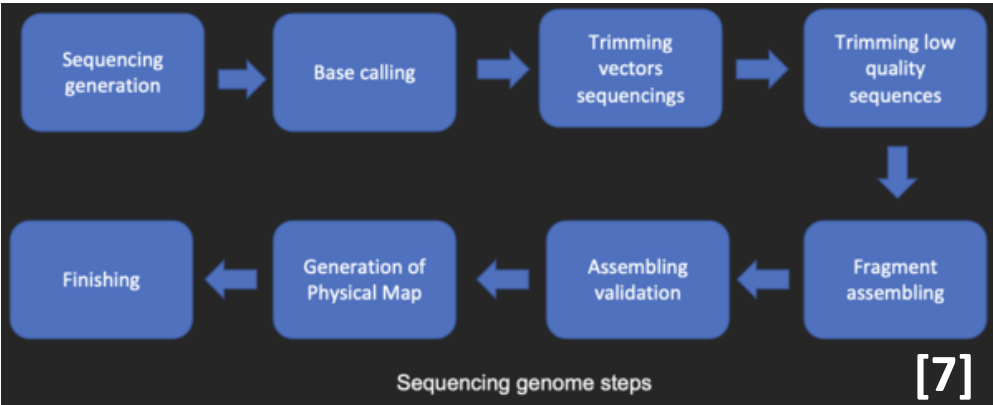
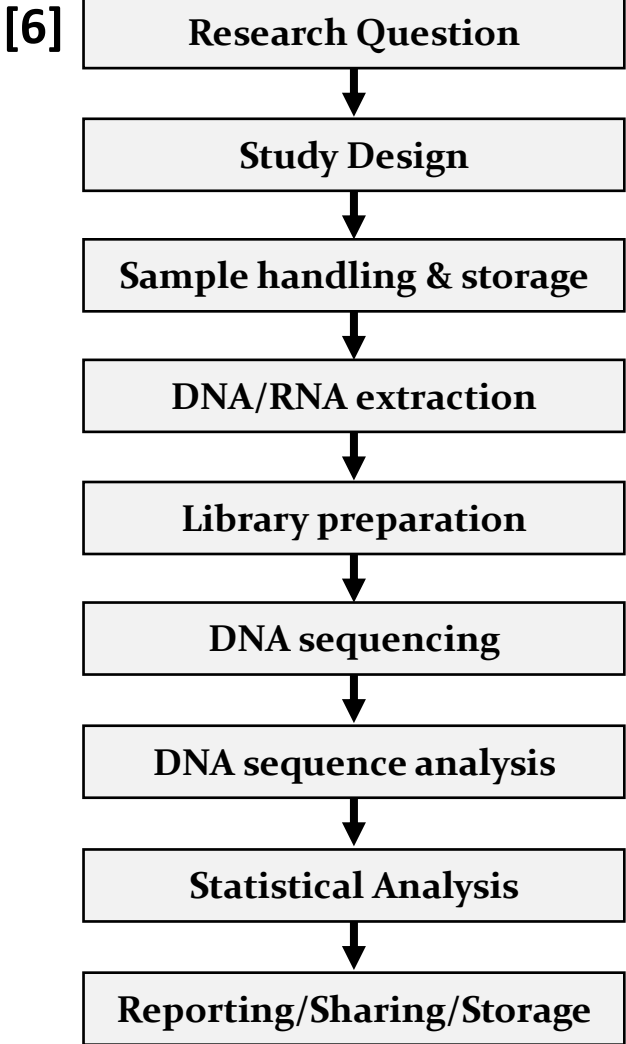
🔧 RNA-Seq

🔧 Microarray

🔧 Phylogeny

🔧 Networks

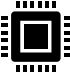
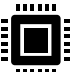
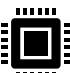
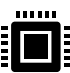
🔧 Omics





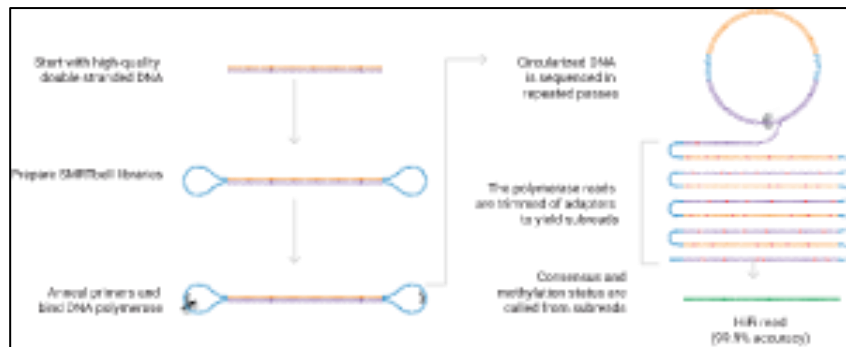
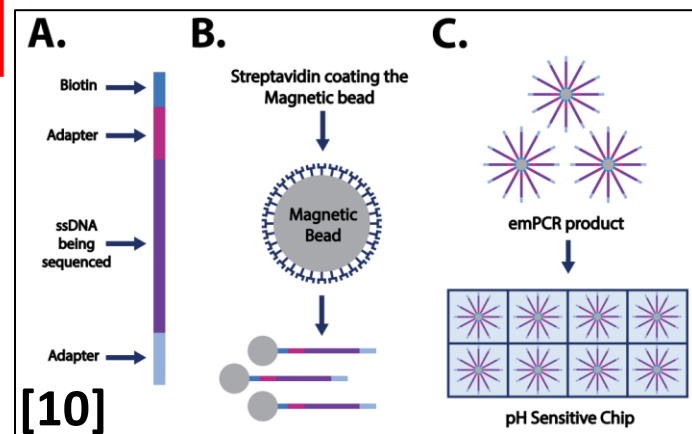
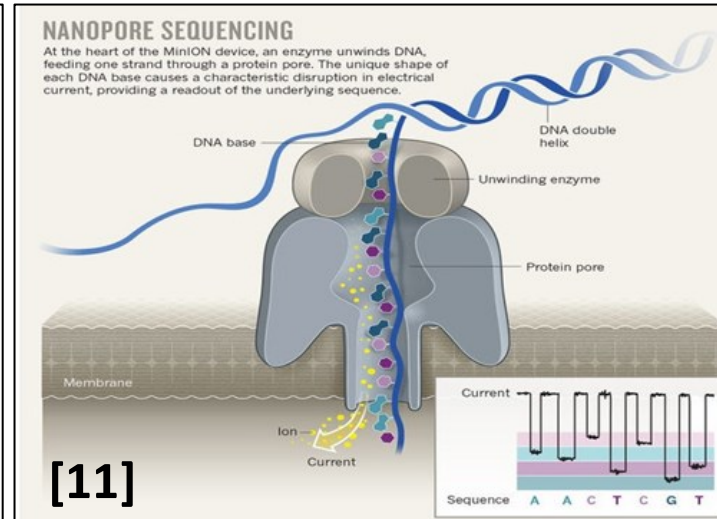
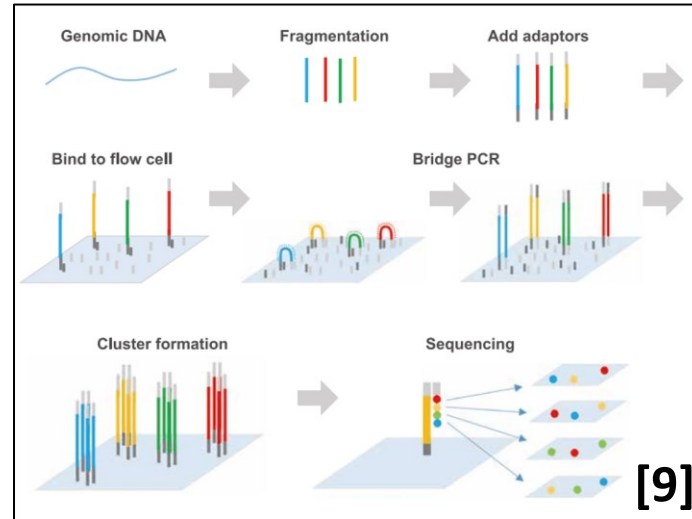
# Machines

 DNA Sequencing [8]:

-  Illumina
-  Ion Torrent
-  Oxford Nanopore
-  Pacific Biosystems

 Mass Spectrometry

 Flow Cytometry



[12]

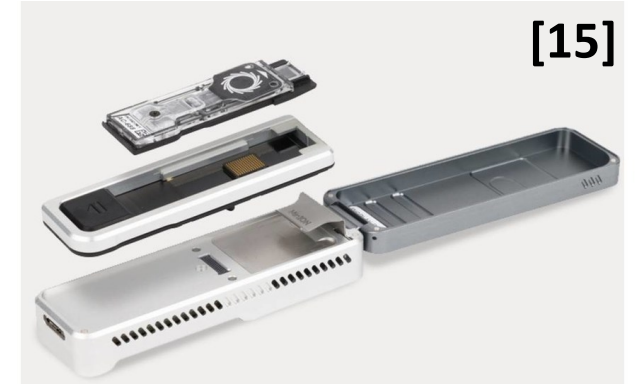
# Machines

## 📡 DNA Sequencing [8]:

- 📡 Illumina
- 📡 Ion Torrent
- 📡 Oxford Nanopore
- 📡 Pacific Biosystems

## 📡 Mass Spectrometry

## 📡 Flow Cytometry



# Machines

 DNA Sequencing [8]:

 Illumina

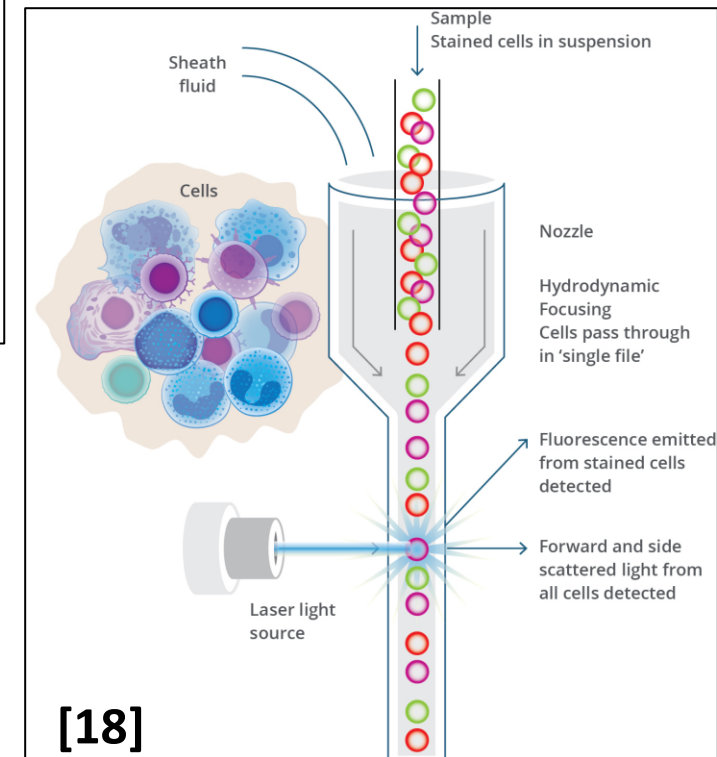
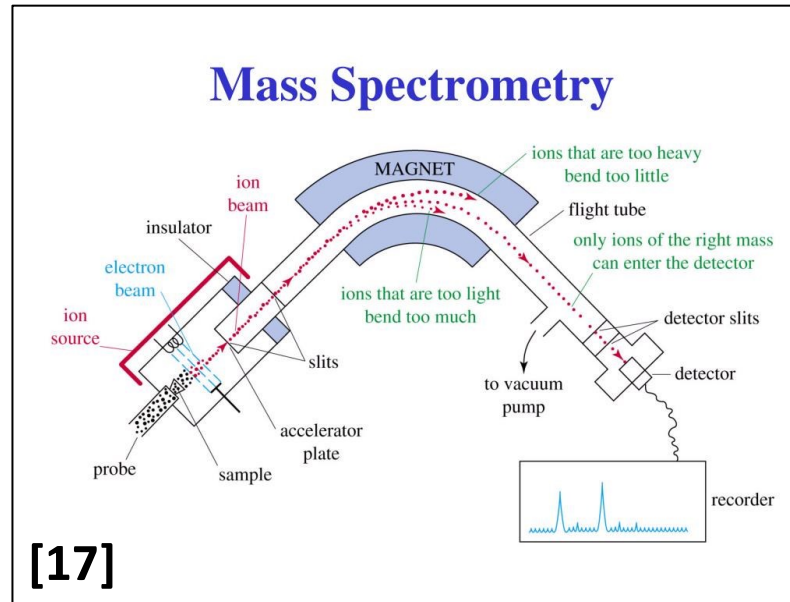
 Ion Torrent

 Oxford Nanopore

 Pacific Biosystems

 **Mass Spectrometry**

 **Flow Cytometry**



# Software



**Core Tools:** BLAST, EMBOSS, Clustal Omega, PROSPECT, Ensembl, Usearch, SAMtools, Velvet, Bowtie, TopHat, Muscle, ProteinTools, etc.



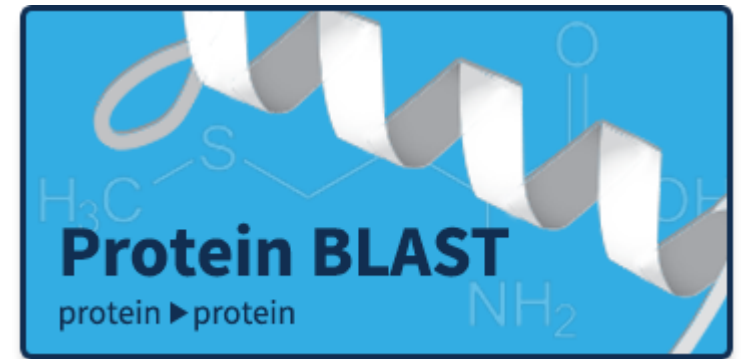
**Underlying Code:** Needleman-Wunch, Borrows-Wheeler Transform, UPGMA, Neighbor-joining, etc.



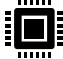
**Software Environments:** Biopython, Bioconductor, BioPerl, etc.

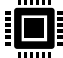


**Pipelines:** Dada2, DNA-seq Analysis, Center for Genomic Epidemiology, Qiime2, Mothur, Workflow Managers, etc.

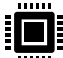


# Software

 **Core Tools:** BLAST, EMBOSS, Clustal Omega, PROSPECT, Ensembl, Usearch, SAMtools, Velvet, Bowtie, TopHat, Muscle, ProteinTools, etc.

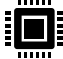
 **Underlying Code:** Needleman-Wunch, Borrows-Wheeler Transform, UPGMA, Neighbor-joining, etc.

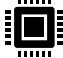
 **Software Environments:** Biopython, Bioconductor, BioPerl, etc.

 **Pipelines:** Dada2, DNA-seq Analysis, Center for Genomic Epidemiology, Qiime2, Mothur, Workflow Managers, etc.

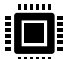


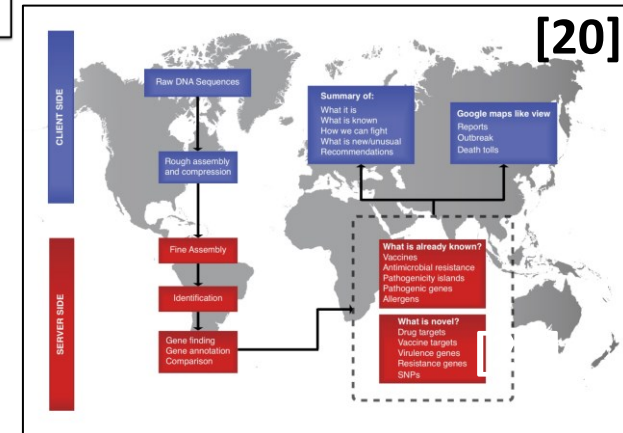
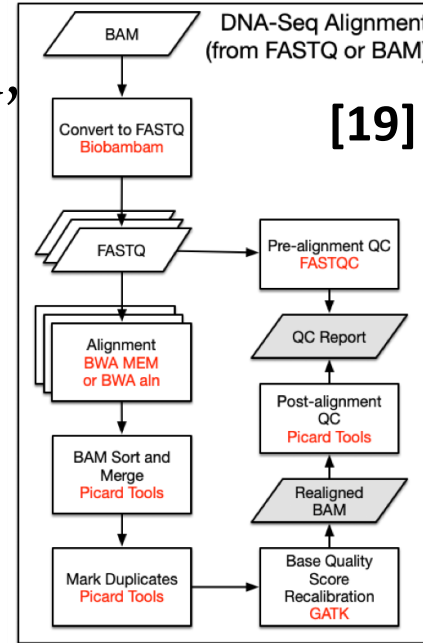
# Software

 **Core Tools:** BLAST, EMBOSS, Clustal Omega, PROSPECT, Ensembl, Usearch, SAMtools, Velvet, Bowtie, TopHat, Muscle, ProteinTools, etc.

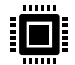
 **Underlying Code:** Needleman-Wunch, Borrows-Wheeler Transform, UPGMA, Neighbor-joining, etc.

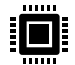
 **Software Environments:** Biopython, Bioconductor, BioPerl, etc.

 **Pipelines:** Dada2, DNA-seq Analysis, Center for Genomic Epidemiology, Qiime2, Mothur, Workflow Managers, etc.

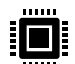


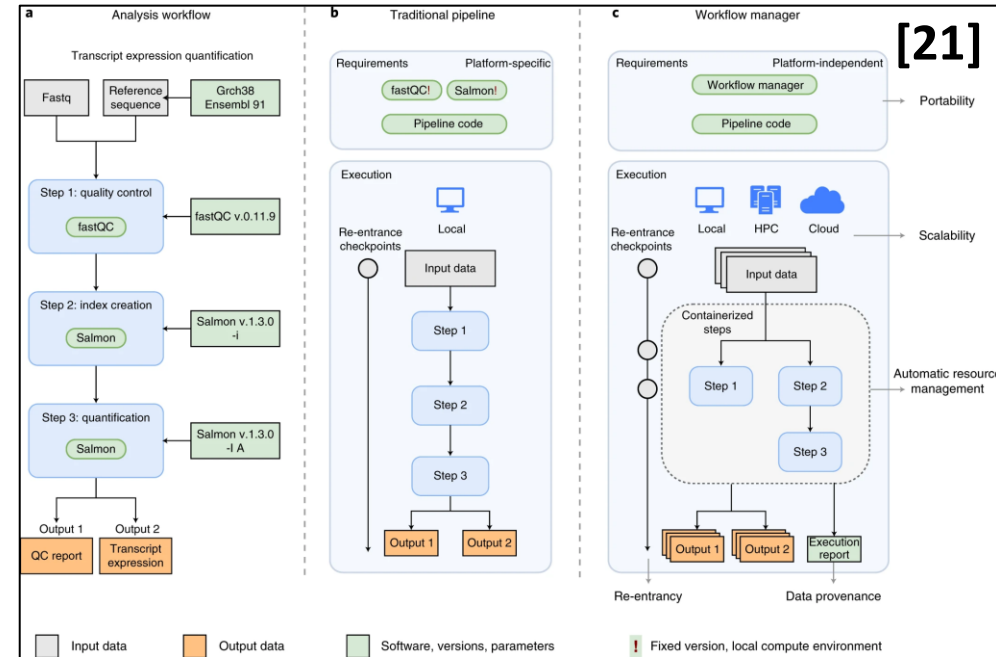
# Software

 **Core Tools:** BLAST, EMBOSS, Clustal Omega, PROSPECT, Ensembl, Usearch, SAMtools, Velvet, Bowtie, TopHat, Muscle, ProteinTools, etc.

 **Underlying Code:** Needleman-Wunch, Borrows-Wheeler Transform, UPGMA, Neighbor-joining, etc.

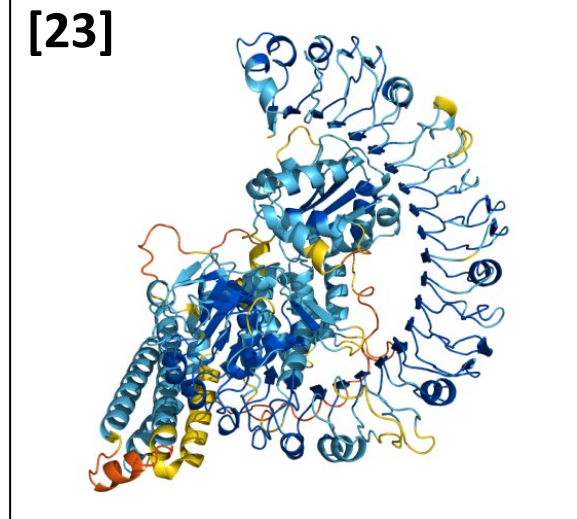
 **Software Environments:** Biopython, Bioconductor, BioPerl, etc.

 **Pipelines:** Dada2, DNA-seq Analysis, Center for Genomic Epidemiology, Qiime2, Mothur, **Workflow Managers, etc.**



# Databases

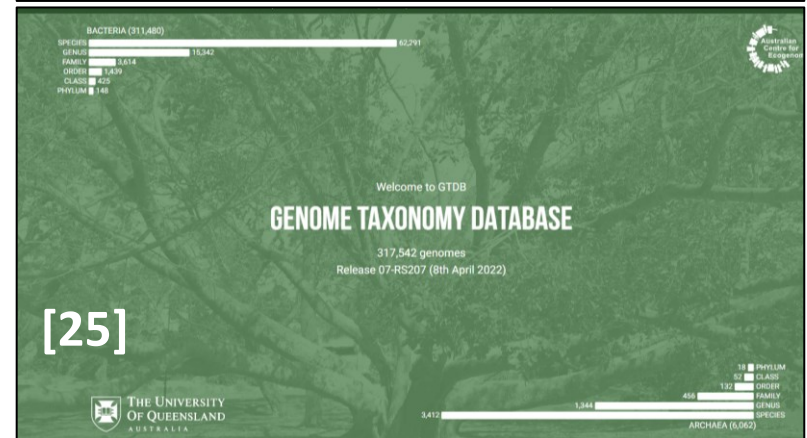
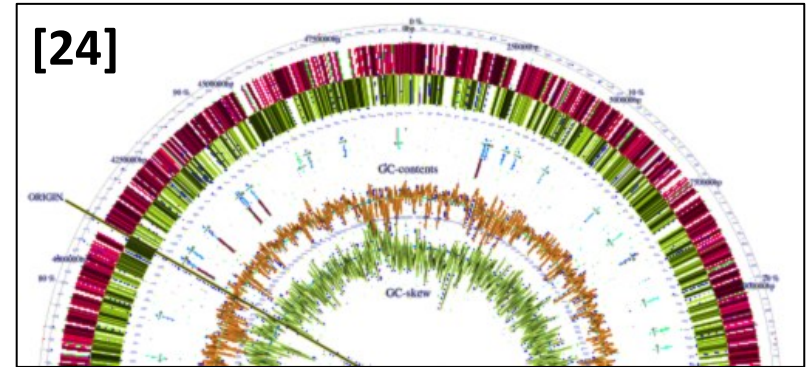
- ❑ **Entities:** NCBI, EMBL, Japanese Institute of Genetics
- ❑ **DNA:** GenBank, RefSeq
- ❑ **Proteins:** UniProt, Alphafold
- ❑ **Pathways:** KEGG
- ❑ **Taxonomy:** NCBI Taxonomy, Catalogue of Life
- ❑ **Interesting:** rrnDB, GTDB, silva



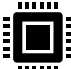
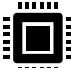
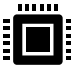
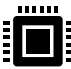


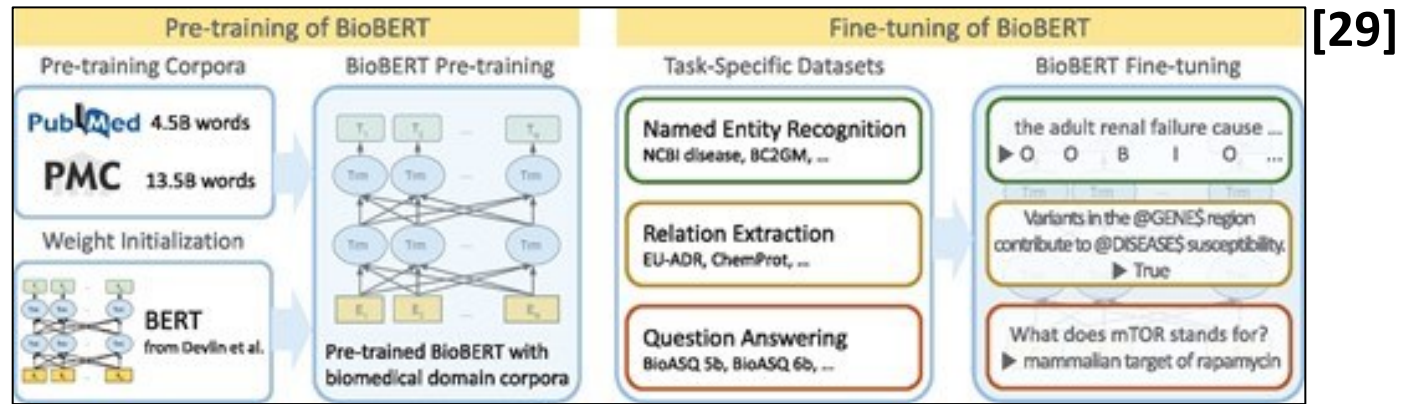
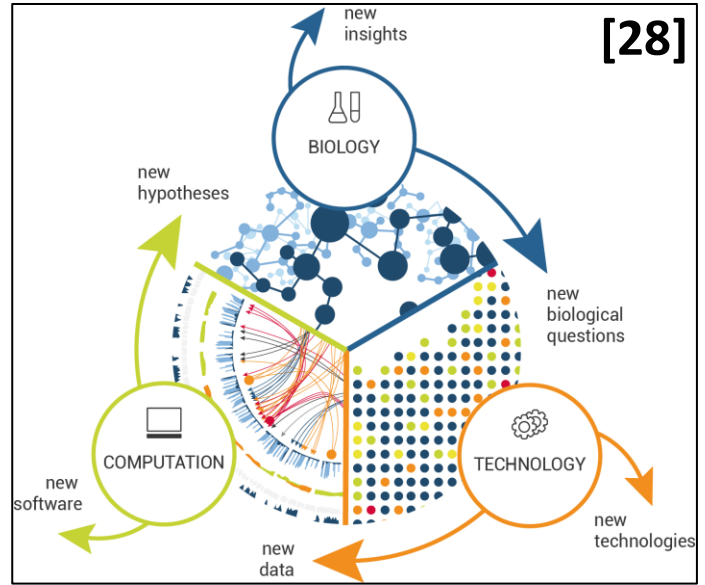
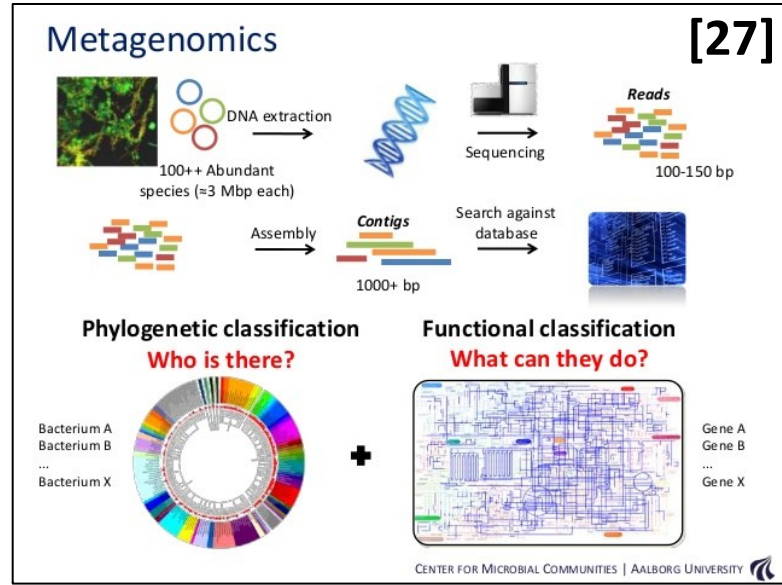
# Databases

- Entities: NCBI, EMBL, Japanese Institute of Genetics
- DNA: GenBank, RefSeq
- Proteins: UniProt, Alphafold
- Pathways: KEGG
- Taxonomy: NCBI Taxonomy, Catalogue of Life
- Interesting: rrnDB, GTDB, silva**

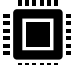
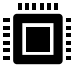
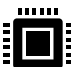
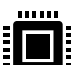


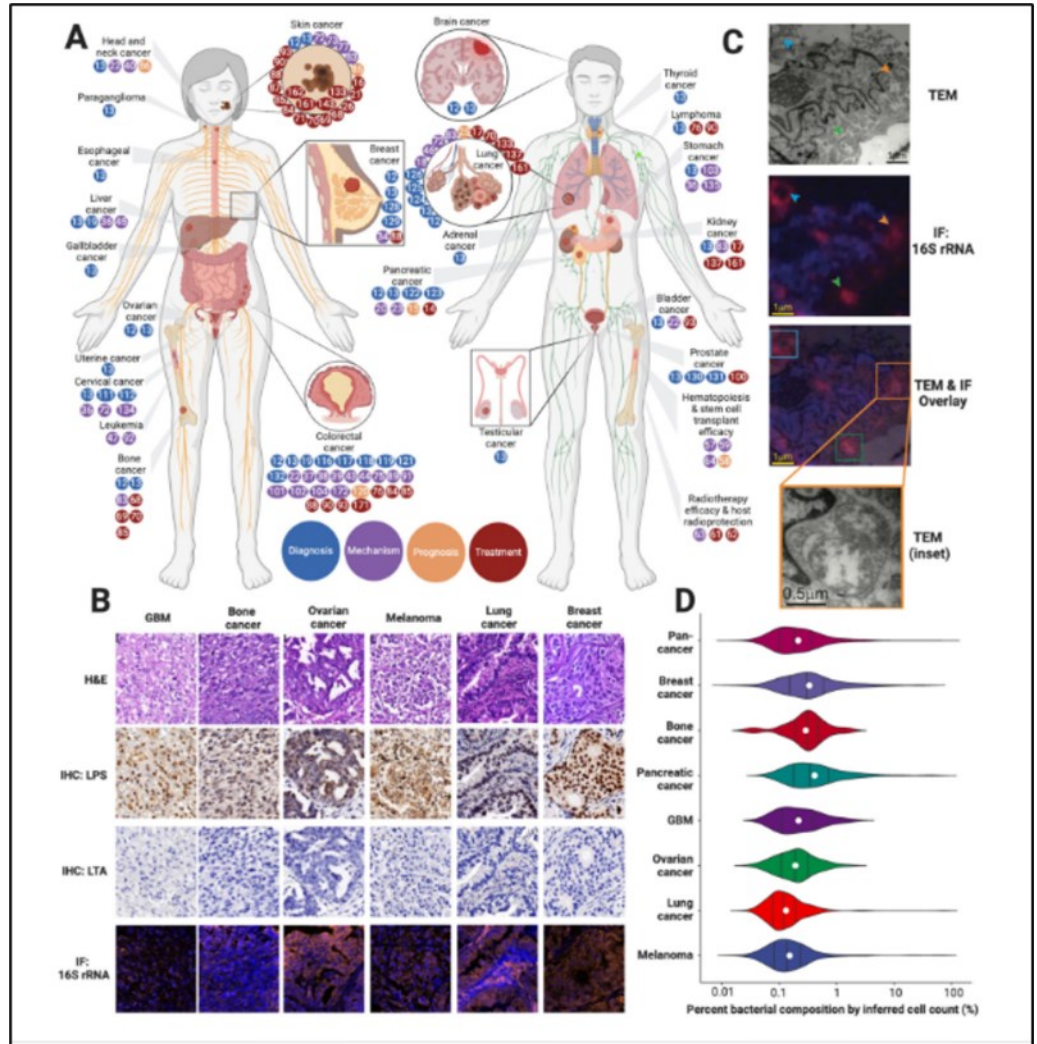
# Example Applications

-  Metagenomics
-  Systems Biology
-  Literature Analysis
-  Microbiome



# Example Applications

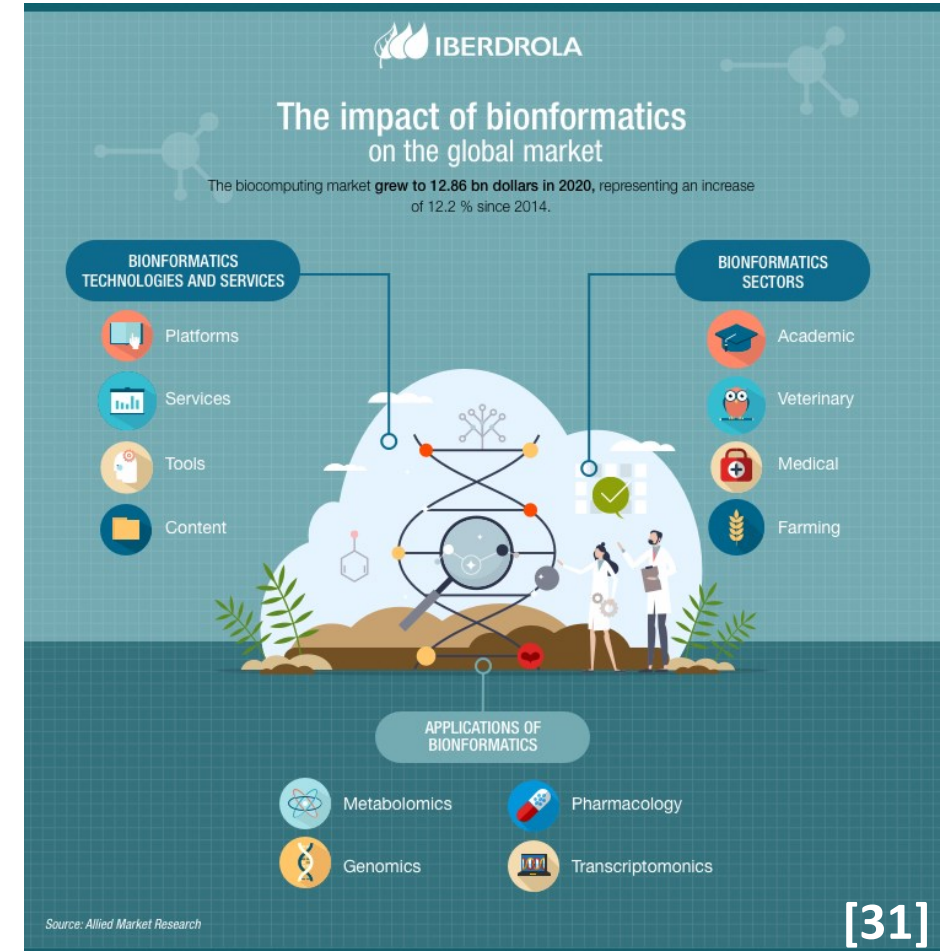
-  Metagenomics
-  Systems Biology
-  Literature Analysis
-  **Microbiome**



[30]

# The Future

- Increasing market share
- Use of robotics, machine learning, deep learning, AI, etc. [32]
- More data, more integration
- The future is long? [33]



# Resources

- ❑ Algorithms for DNA Sequencing
- ❑ Joint Genome Institute
  - ❑ NeLLi 2021 [35]
  - ❑ VEGA Symposium
  - ❑ Annual Genomics of Energy & Environment Meeting
- ❑ Center for Genomic Epidemiology (TDU)
  - ❑ Metagenomics applied to surveillance of pathogens and antimicrobial resistance
  - ❑ Whole genome sequencing of bacterial genomes-tools and applications
- ❑ Bioinformatics.ca [36]

A video thumbnail for a presentation titled "Algorithms for DNA Sequencing" by Ben Langmead. The thumbnail shows a man speaking in front of a screen displaying the text "First generation DNA sequencing" and images of journal covers from "nature" and "SCIENCE".

**Algorithms for DNA Sequencing**

Ben Langmead

[34]

# References 1

- [1] <https://pubmed.ncbi.nlm.nih.gov/30084940/>
- [2] <https://collegedunia.com/exams/bioinformatics-biology-articleid-1454>
- [3] [https://sites.google.com/site/obenscience7e/\\_/rsrc/1490190066963/unit-7/from-dna-to-protein/dna%20to%20rna%20to%20protein.png](https://sites.google.com/site/obenscience7e/_/rsrc/1490190066963/unit-7/from-dna-to-protein/dna%20to%20rna%20to%20protein.png)
- [4] <https://www.ebi.ac.uk/ols/ontologies/edam>
- [5] <https://www.formbio.com/blog/your-essential-guide-different-file-formats-bioinformatics>
- [6] <https://www.coursera.org/learn/wgs-bacteria/>
- [7] <https://en.wikipedia.org/wiki/Bioinformatics>
- [8] <https://genohub.com/ngs-instrument-guide/>
- [9] <https://blogs.iu.edu/ncgas/files/2021/09/Picture1.png>
- [10] <https://i0.wp.com/apollo-institute.org/wp-content/uploads/2021/11/IonTSeq-Figure-1-01-01.png?resize=1024%2C636&ssl=1>
- [11] <https://www.medgadget.com/wp-content/uploads/2020/04/nanopore.png>
- [12] <https://www.pacb.com/technology/hifi-sequencing/>



# References 2

- [13] [https://img.medicalexpo.com/images\\_me/photo-g/83632-12919744.jpg](https://img.medicalexpo.com/images_me/photo-g/83632-12919744.jpg)
- [14] <http://genearrays.com/wp-content/uploads/2014/03/pgm.jpg>
- [15] [https://mma.prnewswire.com/media/834434/Oxford\\_Nanopore\\_Flongle.jpg?p=facebook](https://mma.prnewswire.com/media/834434/Oxford_Nanopore_Flongle.jpg?p=facebook)
- [16] <https://www.pacb.com/revio/>
- [17] <https://image.slideserve.com/464232/mass-spectrometry-l.jpg>
- [18] <https://www.streck.com/wp-content/uploads/2020/11/flow-cytometry-light-scatter-illustration.jpg>
- [19] [https://docs.gdc.cancer.gov/Data/Bioinformatics\\_Pipelines/DNA\\_Seq\\_Variant\\_Calling\\_Pipeline/](https://docs.gdc.cancer.gov/Data/Bioinformatics_Pipelines/DNA_Seq_Variant_Calling_Pipeline/)
- [20] <https://www.genomicepidemiology.org/>
- [21] <https://www.nature.com/articles/s41592-021-01254-9/figures/1>
- [22] By Anton Nekrutenko, <http://galaxyproject.org/>, Fair use, <https://en.wikipedia.org/w/index.php?curid=32778899>
- [23] <https://alphafold.ebi.ac.uk/>
- [24] <https://rrndb.umms.med.umich.edu/>



# References 3

- [25] <https://gtdb.ecogenomic.org/>
- [26] <https://www.arb-silva.de/>
- [27] <http://myriverside.sd43.bc.ca/gracynk2015/files/2016/11/MEtagenomics-1-245mzc9.jpg>
- [28] <https://www.omicscouts.com/media/files/blockcontent/2016-09/SystemsBiology.jpg>
- [29] <https://academic.oup.com/bioinformatics/article/36/4/1234/5566506?login=false>
- [30] <https://pubmed.ncbi.nlm.nih.gov/33766858/>
- [31] <https://www.iberdrola.com/innovation/bioinformatics>
- [32] <https://www.linkedin.com/pulse/future-bioinformatics-trends-predictions-venkatesh-chellappa/>
- [33] <https://www.youtube.com/watch?v=KVz6UtNaWbE&list=PL2mpR0RYFQsBiCWVJSvVAO3OJ2t7DzoHA&index=54>
- [34] <https://www.youtube.com/playlist?list=PL2mpR0RYFQsBiCWVJSvVAO3OJ2t7DzoHA>
- [35] <https://www.youtube.com/playlist?list=PLkxZMDuKlaKs80mdPGKRI1r7DKII1O8L1>
- [36] <https://bioinformaticsdotca.github.io/>





# Acknowledgements



The DaCCoTA is supported by the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729.

For the labs that use the Biostatistics, Epidemiology, and Research Design Core in any way, including this Module, please acknowledge us for publications: *"Research reported in this publication was supported by DaCCoTA (the National Institute of General Medical Sciences of the National Institutes of Health under Award Number U54GM128729)"*

**DaCCoTA**  
DAKOTA COMMUNITY COLLABORATIVE  
ON TRANSLATIONAL ACTIVITY

