

# ComBINEs

Computational Biology and  
Bioinformatics Network Stuttgart

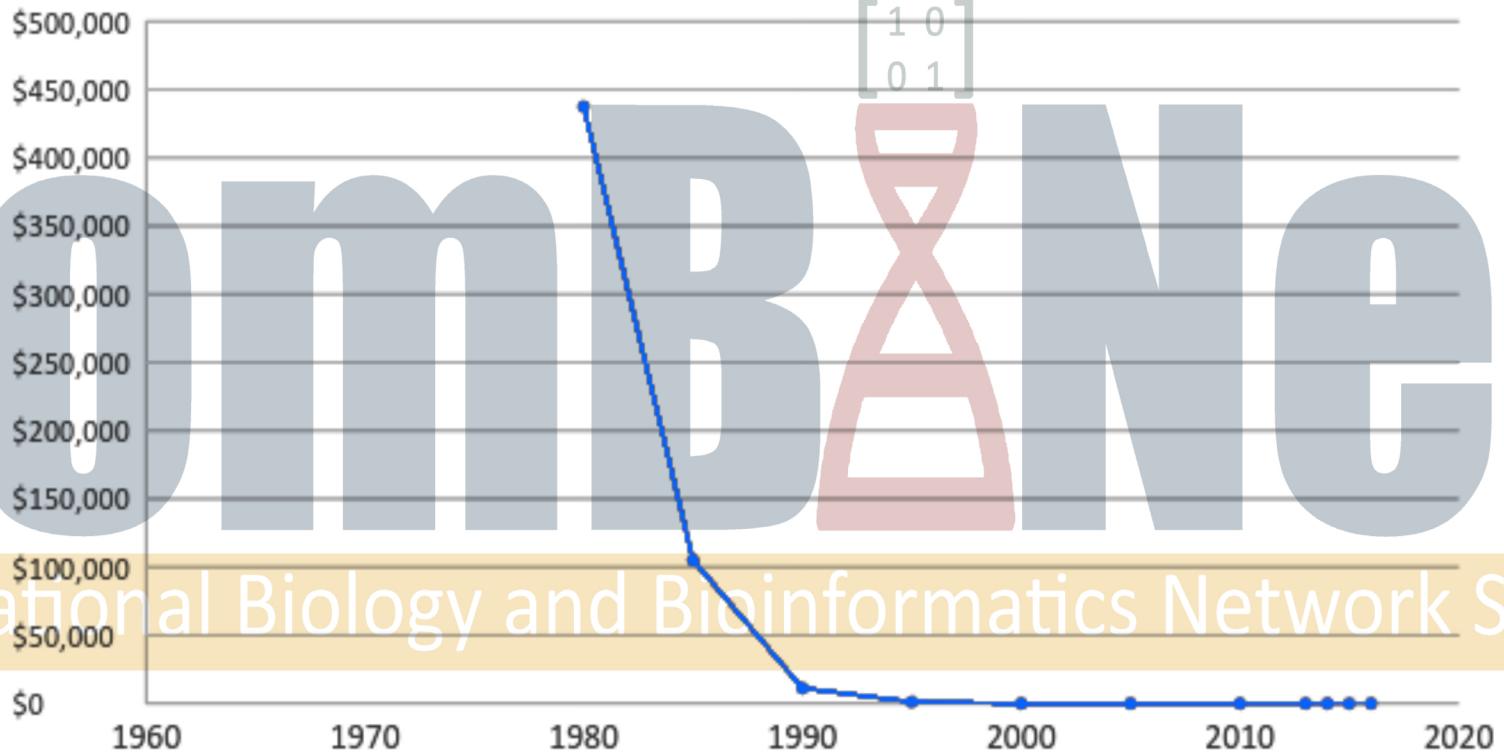


Big Data in Biology & Medicine @ 06.06.2019  
Dr. Nandor Poka

Computational Biology and Bioinformatics Network Stuttgart

# Big Data

Average Cost Per Gigabyte

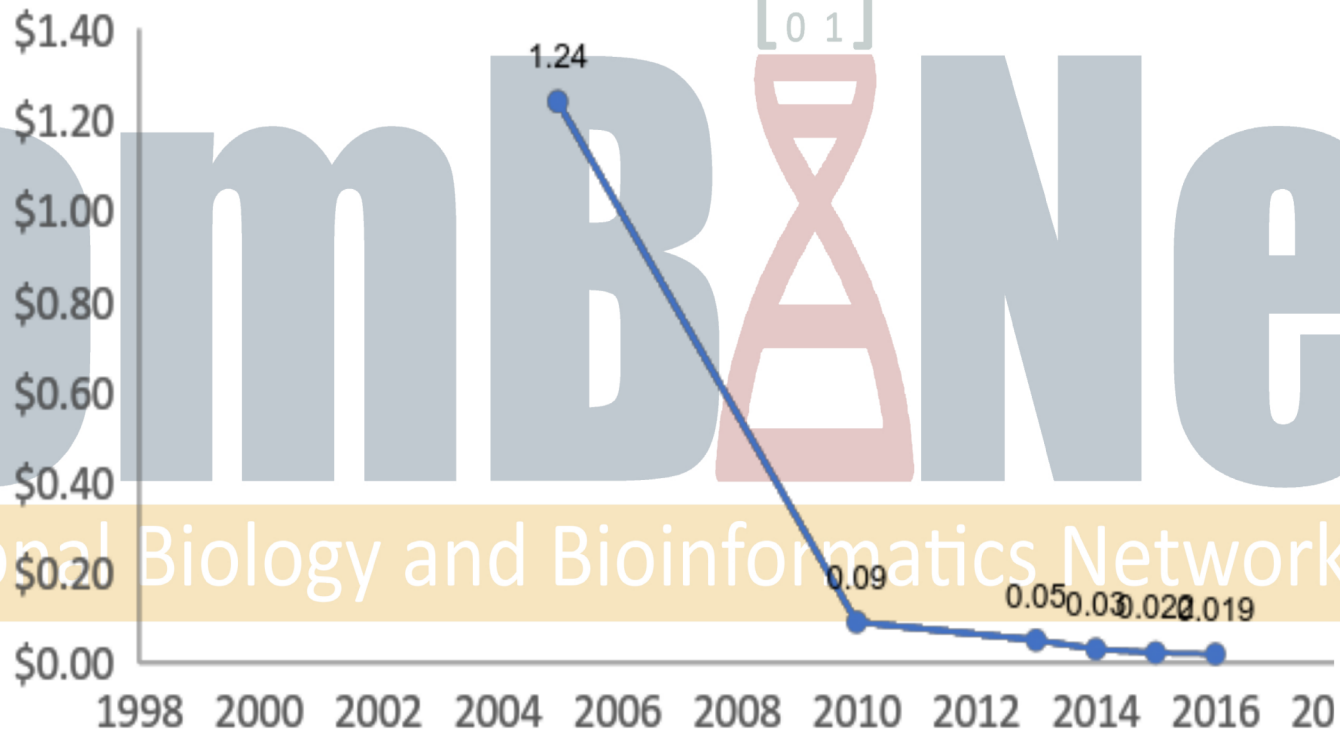


Computational Biology and Bioinformatics Network Stuttgart

# ComBBNES

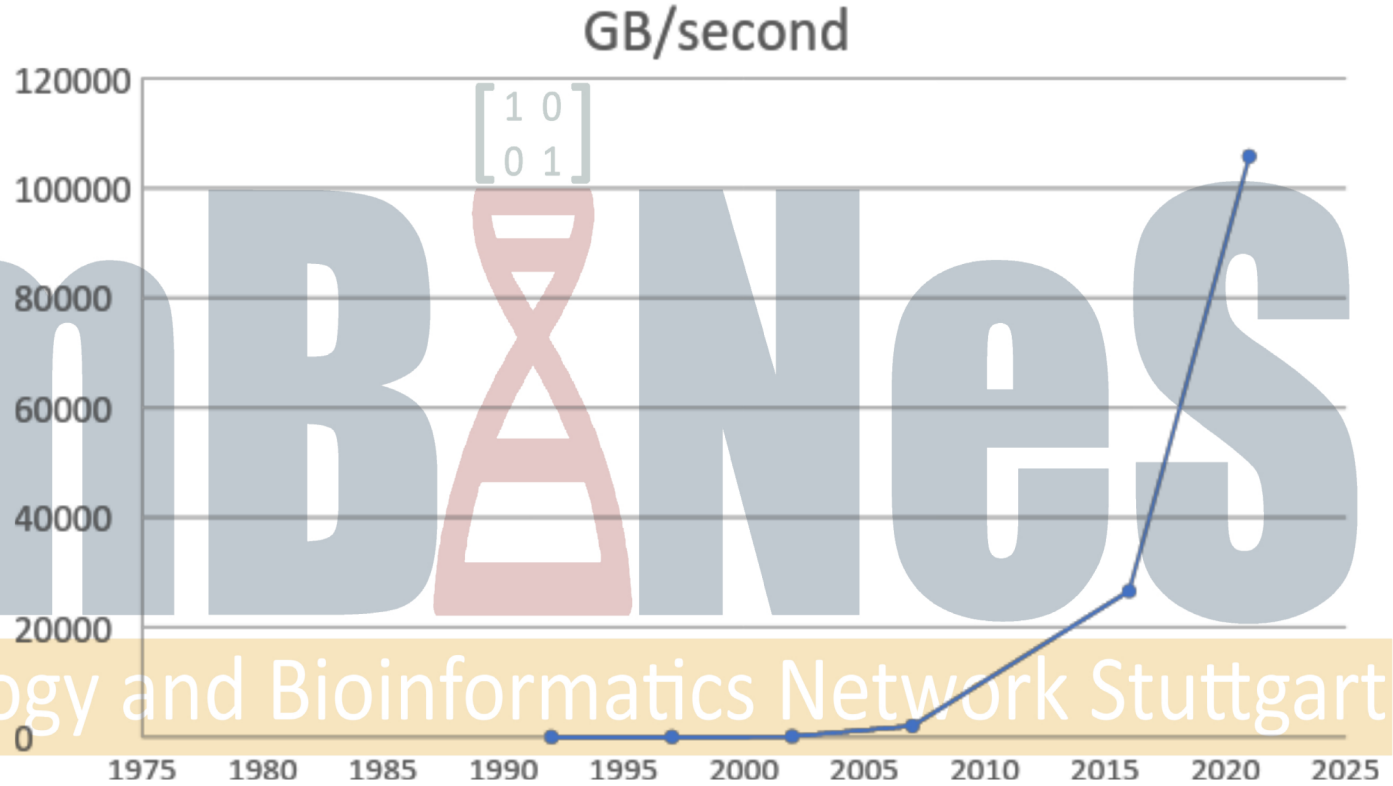


Computational Biology and Bioinformatics Network Stuttgart



**Year Global Internet Traffic**

1992	100 GB per day
1997	100 GB per hour
2002	100 GB per second
2007	2,000 GB per second
2016	26,600 GB per second
2021	105,800 GB per second



# But what is Big Data?

- Like old data, just a lot more... too much to deal with traditional methods
- System logs of various kinds
- Financial transactions
- Online purchases
- Etc.

Computational Biology and Bioinformatics Network Stuttgart

# “New age” Big Data

- Social media – likes, posts, contacts
- User location data – navigation searches, phone location data
- Community knowledge – recommendations, product ratings
- Personal (non-professional) health data – fitness trackers

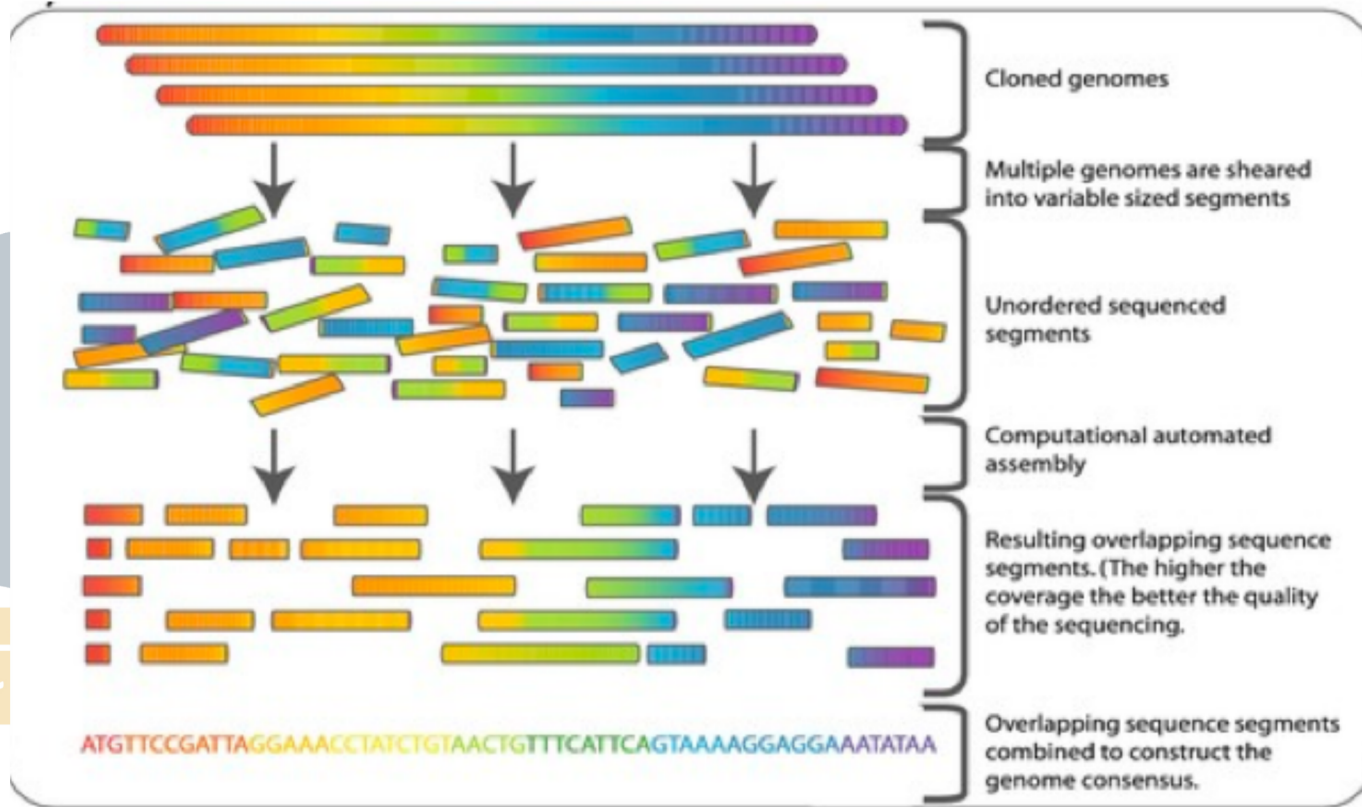
Computational Biology and Bioinformatics Network Stuttgart

# What sources of data we have in Biology and Medicine?

- Sequencing data
- Genomics / proteomics (and other “omics”) data
- Gene / protein interaction networks
- Medical images
- Electronic personal medical data
- Wearable fitness data

Computational Biology and Bioinformatics Network Stuttgart

# Sequencing



CS

Computat

CS

Stuttgart

Source:

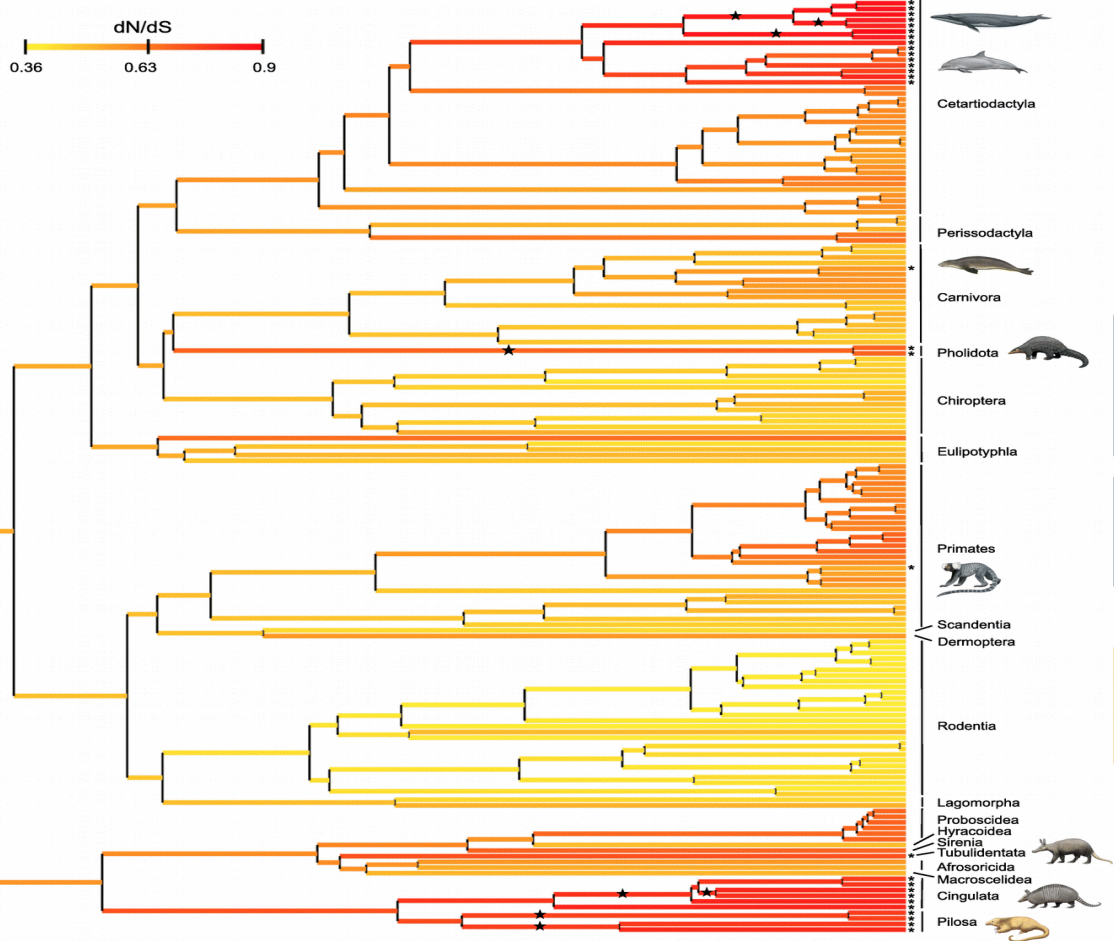
Computational Biology Methods and Their Application to the Comparative Genomics of Endocellular Symbiotic Bacteria of Insects. Biol. Procedures Online (2009). Commins, J., Toft, C., Fares, M. A.



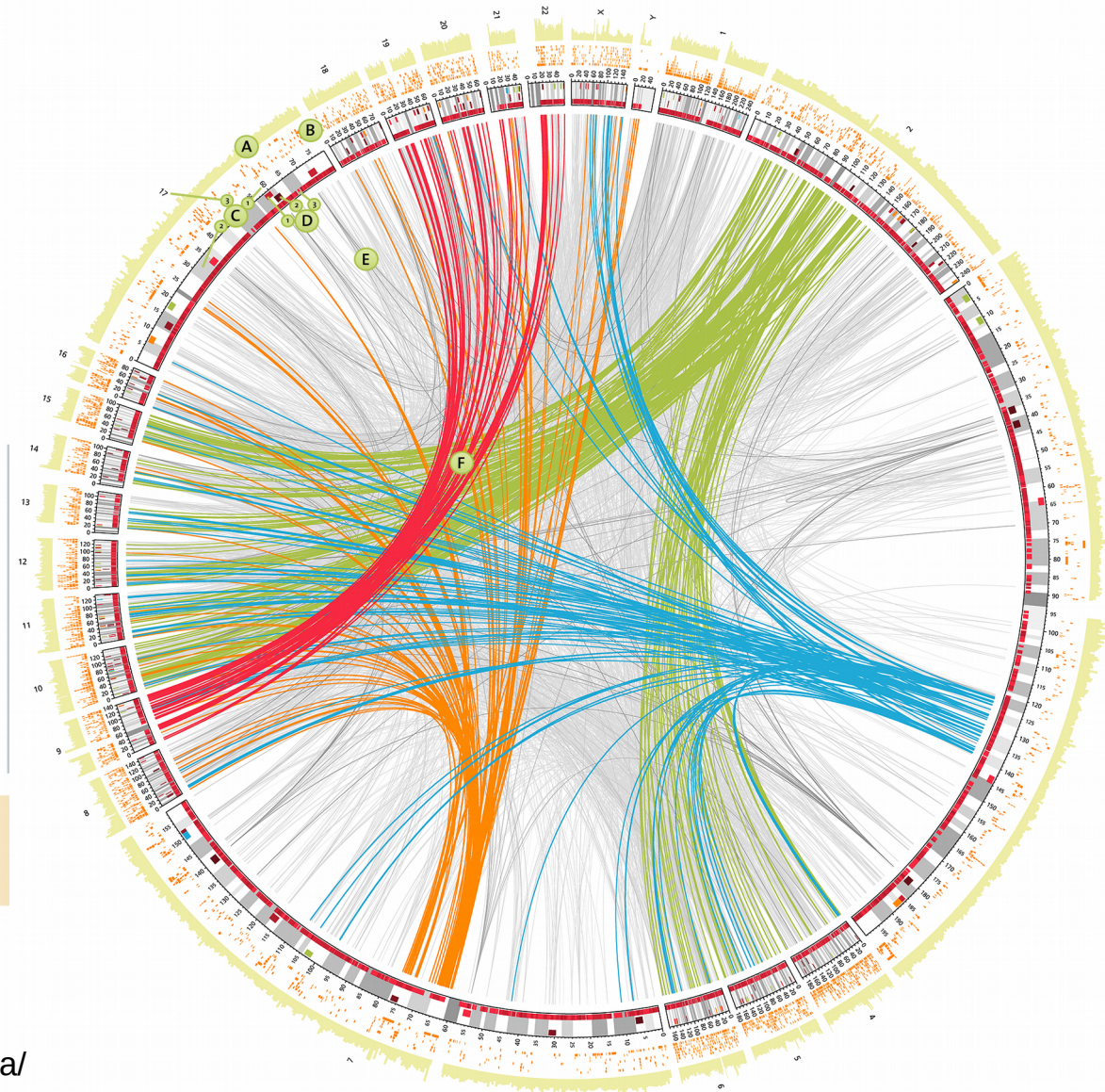
# Genomics and other “omics”

COIN  
 Computational Bio

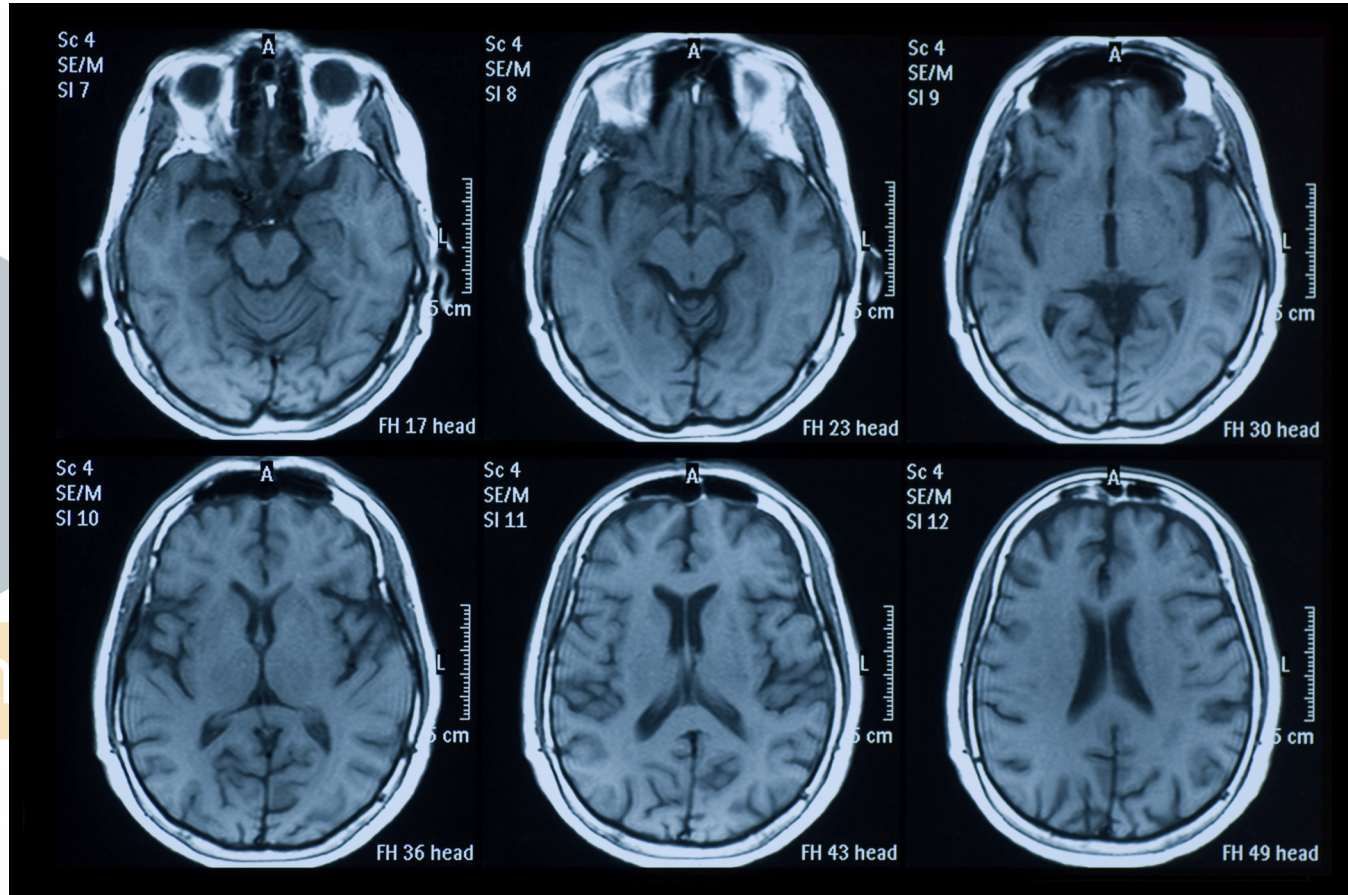
IS  
 Stuttgart



Source:  
<https://doi.org/10.1186/s12862-019-1359-6>



# Medical imaging



Co

Computation

S

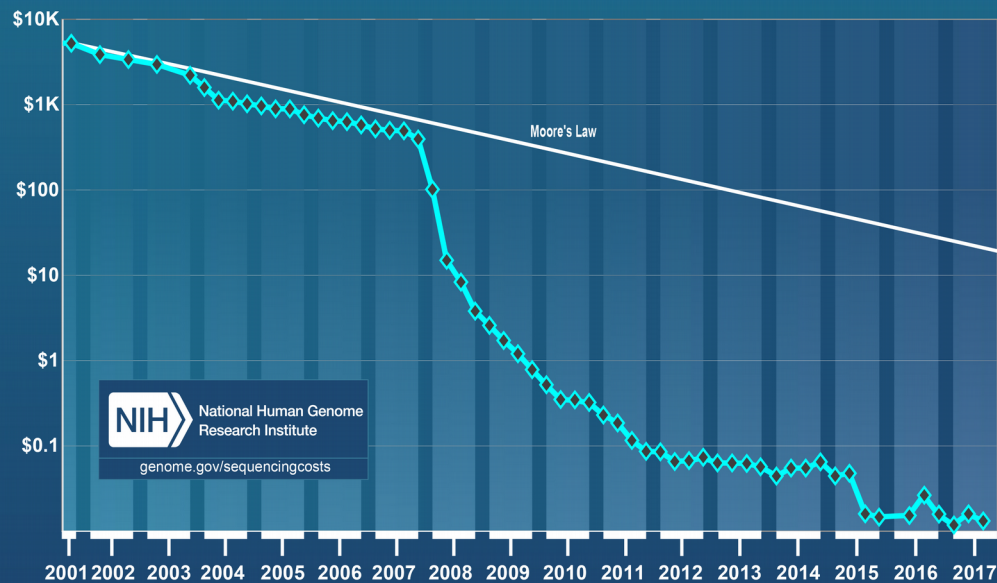
Stuttgart

# What type of data?

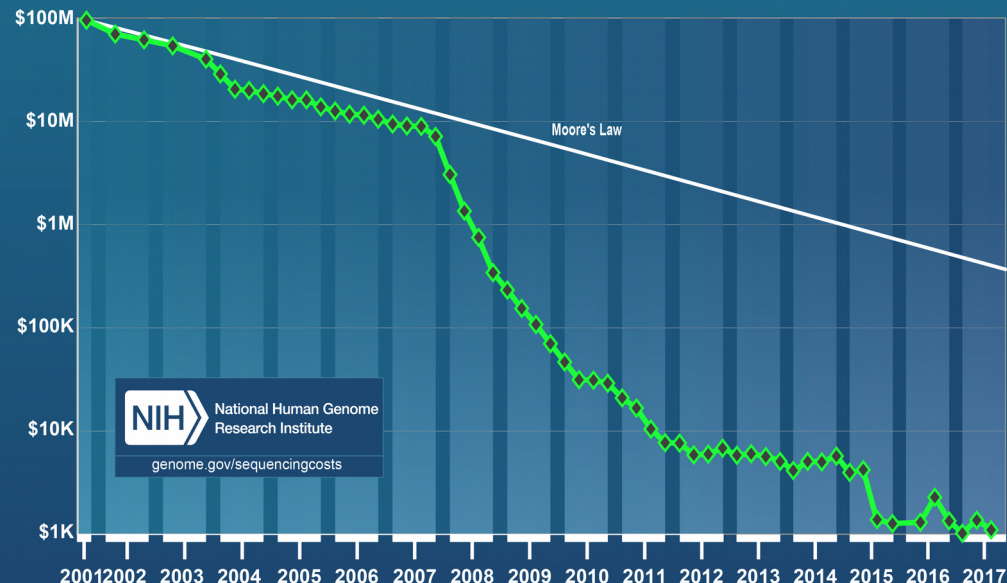
- Sequences → plain text, often compressed
- “Omics” data → heterogenous, but typically includes sequence data (often with added metadata / markup), tables / databases, graphs, 3D molecule models
- Medical images → high-res often multi-layered images / videos
- Medical records → can be anything
- Personal fitness data → varies, but somewhat unified → eg. Google Fit data format / types

# How is this Big Data?

## Cost per Raw Megabase of DNA Sequence



## Cost per Genome



- The human diploid genome is approx 6.6Gb (GigaBase) → in plain ASCII it's 6.6GB (GigaByte)
- Typical sequencing is done with 100 base long reads and 40X coverage →  $\sim 2.64 \cdot 10^9$  reads → 264GB raw sequence
- $\sim 2x$  more with quality data → 0,5TB+
- Some organisms have genome size of 100-200Gb+ genome sizes

Computational Biology and Bioinformatics Network Stuttgart

```
@SEQ_ID
```

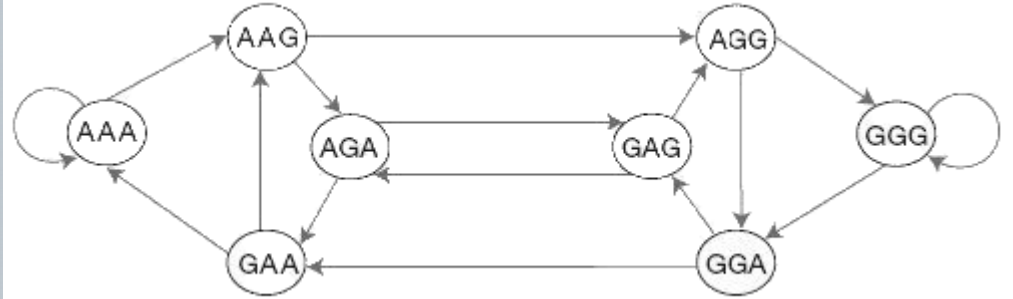
```
GATTTGGGGTTCAAAGCAGTATCGATCA
```

```
+
```

```
!''*(((((***+))%%#+))(%%%).1
```

# For assembly

- de Bruijn graph with 66 Mill. Vertices
- Min. same amount of directed edges



# “Omics” data

- Sample from an gene expression experiment data
- 47191 rows by 405 columns = ~19,1M data points

ID_REF	1NL.AVG_Signal	1NL.Detection Pval	1TU.AVG_Signal	1TU.Detection Pval	2NL.AVG_Signal
ILMN_1762337	61.1	0.24156	59.8	0.31558	64.3
ILMN_2055271	63.3	0.16364	68.8	0.07662	67.4
ILMN_1736007	54.7	0.51818	61.6	0.24545	55.5



# Medical data

- MRI scans → ~3Mpx per image → ~1.7MB as lossless PNG
- 20 images per patient → 24MB / Scan
- In Germany, in 2015, 136.2 MRI scans / 1.000 resident
- $(82.175\text{M} / 1000) * 136.2 = \sim 11.2\text{M Scans} \rightarrow \sim 1,14\text{PB}$   
(Petabyte) of data.

# What can be achieved bio / medicine side?

- “full stack” organism analysis
- Rapid pathogen detection – in plants / animals as well
- Personalized medicine – Cas9 / CRISPR → “on-site” gene repairing
- Rapid drug development

Computational Biology and Bioinformatics Network Stuttgart

# Challenges

- Data privacy
- Cost – manual labour, ingredient and equipment cost
- Quantity vs. quality – data is abundant, but heterogeneous
- Visualization
- Speed(up)
- Free / commercial software

Computational Biology and Bioinformatics Network Stuttgart

# What can be done on IT side?

- Better, faster software / algorithms – cluster / parallel computing
- AI for applied cases – detecting micro fractures on X-rays, detecting tumors on MRI
- Better equipments, data acquirement
- Applied Data analysis / predictions:
  - Live predictions – epilepsy, heart attacks
  - Early onset detection of neurodegenerative diseases (eg. Parkinsons)

# Get your share of the data!

- <https://www.ebi.ac.uk/services>
- <https://www.uniprot.org/>
- <https://www.ncbi.nlm.nih.gov/>
- <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5506686/>

ComBBNet

Computational Biology and Bioinformatics Network Stuttgart

Questions?

