

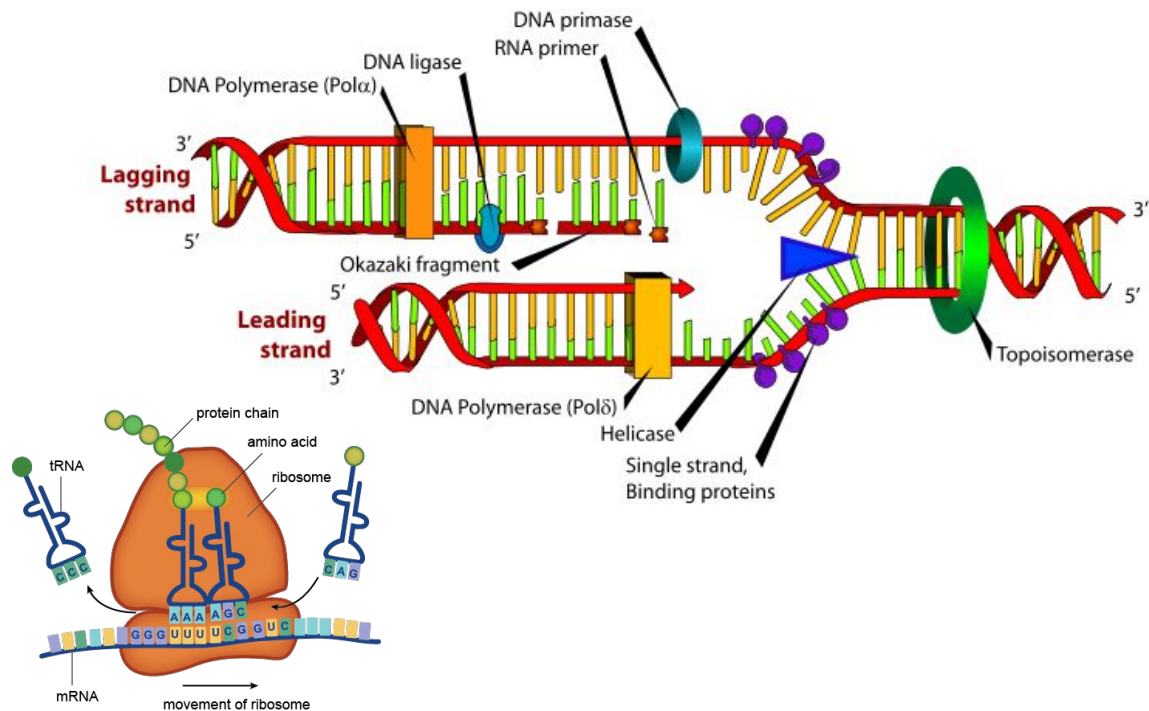
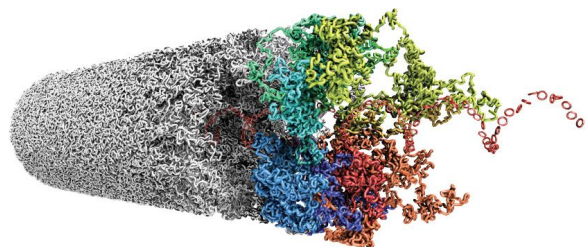
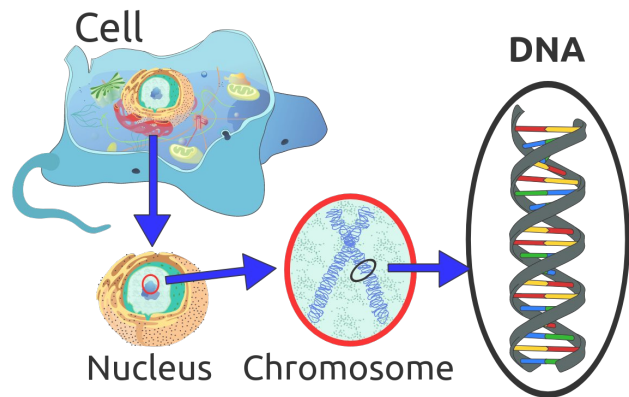
Integrative Gene Selection

-

Intermediate Presentation

Trends in Bioinformatics WS 2018/2019

Biology Recap



Feature Selection



⇒ (Feature Selection) ⇒



- **Short, fat data ($p \gg n$)**

RNAseq data typically covers p -thousands of genes but only few n -hundred samples

- **Incomplete view on cell process**

RNAseq data covers a complete snapshot of the gene activity of a cell; only a single point in time - not more, not less...

- **Gene interactions**

Genes react to changes of partners in the interaction network, interactions have to be taken into account

- **Biological Relevance**

Statistical approaches only concentrate on statistical relevance, ignoring the biological relevance of genes, like disease driver genes, while only showing the respective behaviour of increased gene expression profiles of affected pathway genes

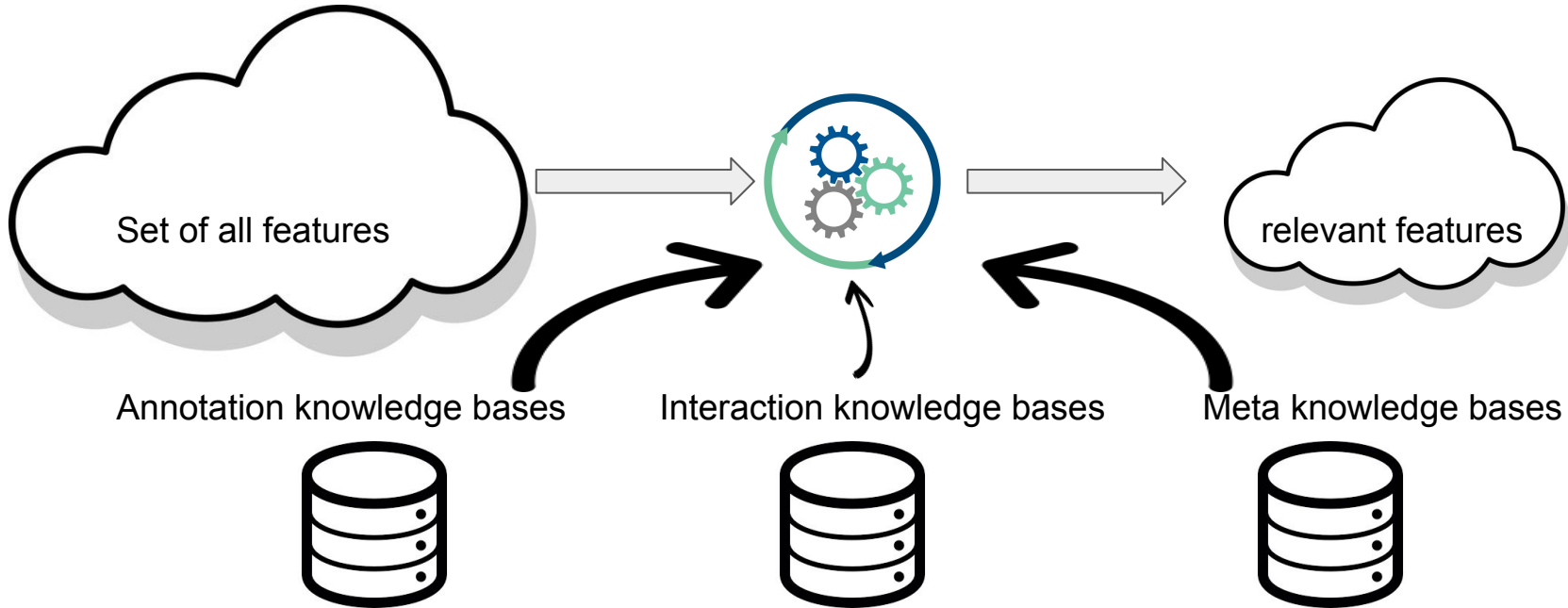
Gene Selection (a.k.a. Feature Selection 4 RNAseq-Data)

- **Goals:**
 - shorten training times
 - avoid the curse of dimensionality (overfitting)
- **Approach:**
 - Use statistical models to select interesting “genes” for future analysis
 - Currently use only a bit of existing knowledge of underlying biological processes

Different approaches exist, depending on where and how they integrate external knowledge!

Integrative Gene Selection

Gene Selection using **not only statistical** methods **but also external knowledge bases!**



There is a **high demand** for
a **widely accepted** and **easily interpretable**
Integrative gene selection solution!

What drives us?

- As by today: plenty problems (though list was shortened)
- Knowledge bases grow constantly, reveal new information
- We do not know much regarding biological processes in humans

What drives us? (ctd.)

- Especially network approaches for gene selection are very promising!
⇒ shown by different studies!
- Compare different integrative approaches
⇒ Framework was built already, missing a network-approach
- Goal: participate in offering researchers a great framework for promising research!

How to solve the problems?

Implement Integrative Gene Selection approach

- using network-based gene extraction similar to Gu et al's [2004] approach
- integrate functional categories or functional modules for specific genes via Gene Ontology (GO)'s REST-API

If time is left, I want to integrate biological pathways as features, also from GO, like Quanz et al [2008] have done.

What will be the next steps?

1. Implement network-based gene selection solution and integrate KB
Functional modules and possibly biological Pathways from GO
2. Integrate solution into existing framework IGEA
3. Evaluate solution; compare approach against existing approaches of IGEA
4. Write paper based on findings of achieved work status

Now is the time
for Questions
and Answers!!!

That's it! I am excited to answer your questions or skip this part and start coding right away!

Thank you for your attention!

Image Sources

- DNA Replication: <http://www.vcbio.science.ru.nl/en/virtuallessons/cellcycle/trans/>
- Chromosome - DNA: <https://www.sciencenews.org/article/dna-chromosome-bundles-cell-division-mitosis>
- Eukaryote: https://da.wikipedia.org/wiki/Fil:Eukaryote_DNA-en.svg
- mRNA to Protein: <https://dumielauxepices.net/wallpaper-3773168>
- Gene Ontology: <http://www.geneontology.org/page/go-citation-policy>
- TCGA: <https://www.genome.gov/17516564/the-cancer-genome-atlas/>
- WikiPathways: <https://en.wikipedia.org/wiki/WikiPathways>
- Database: https://www.flaticon.com/free-icon/database_4426
- Cloud: https://all-free-download.com/free-vector/download/nuage-cloud_116075.html
- Process: <https://positiveenergy.pro/process/>
- Curved Arrow: <https://thenounproject.com/term/curved-arrow/69289/>

Additional Slides, if needed

Attachment

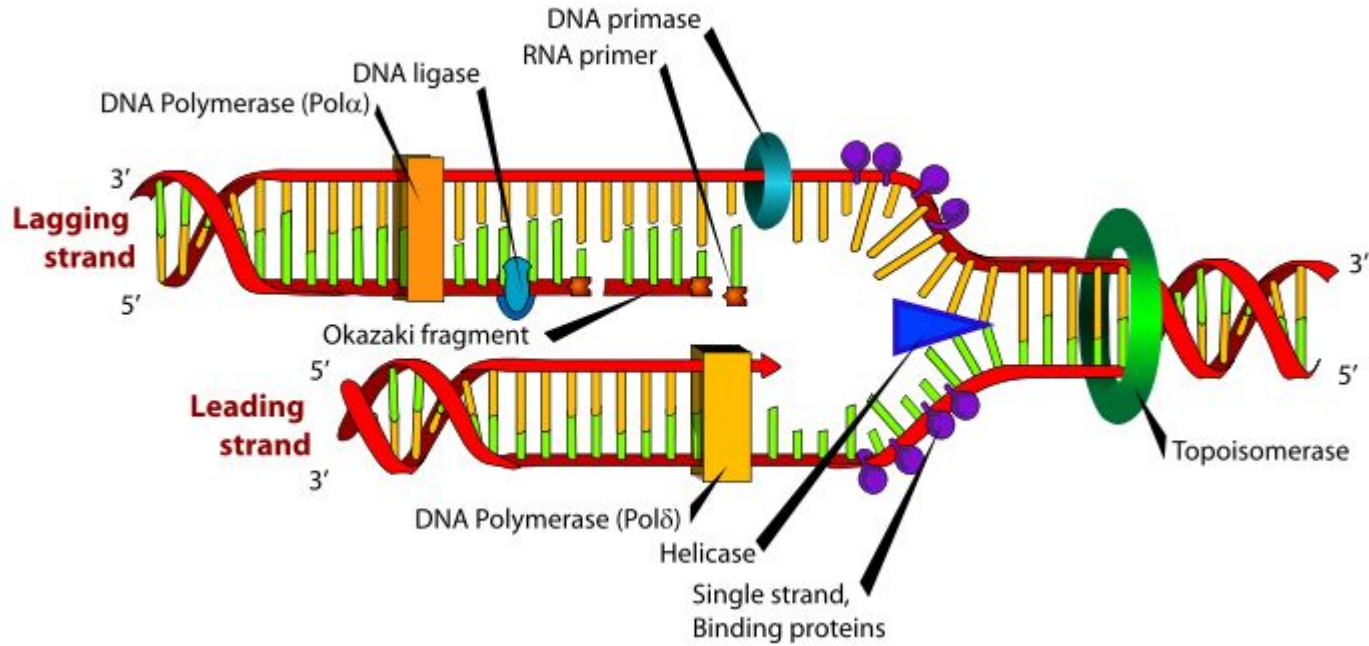
RNA-Seq

RNA-Seq is a recently developed transcriptome profiling technology that utilizes next-generation sequencing platforms (Metzker, 2010; Mardis, 2008). RNA-Seq transcripts are reverse-transcribed into cDNA, and adapters are ligated to each end of the cDNA. Sequencing can be done either unidirectional (single-end sequencing) or bidirectional (paired-end sequencing) and then aligned to a reference genome database or assembled to obtain *de novo* transcripts, proving a genome-wide expression profile (Wang et al., 2009). RNA-Seq offers many advantages over microarray technology. Unlike microarray technology, which depends on already known genes, RNA-Seq is not dependent on existing genome data and can screen novel transcript and analyze transcript structure, including single base-pair resolution and exonic boundaries, which is very valuable while investigating SNPs, thus making it useful for genotyping and linkage analysis (Wang et al., 2009). The advantages of RNA-Seq and its application in studying nervous system and the challenges associated with the technology are summarized in a previous publication (Kadakkuzha and Puthanveetil, 2013).

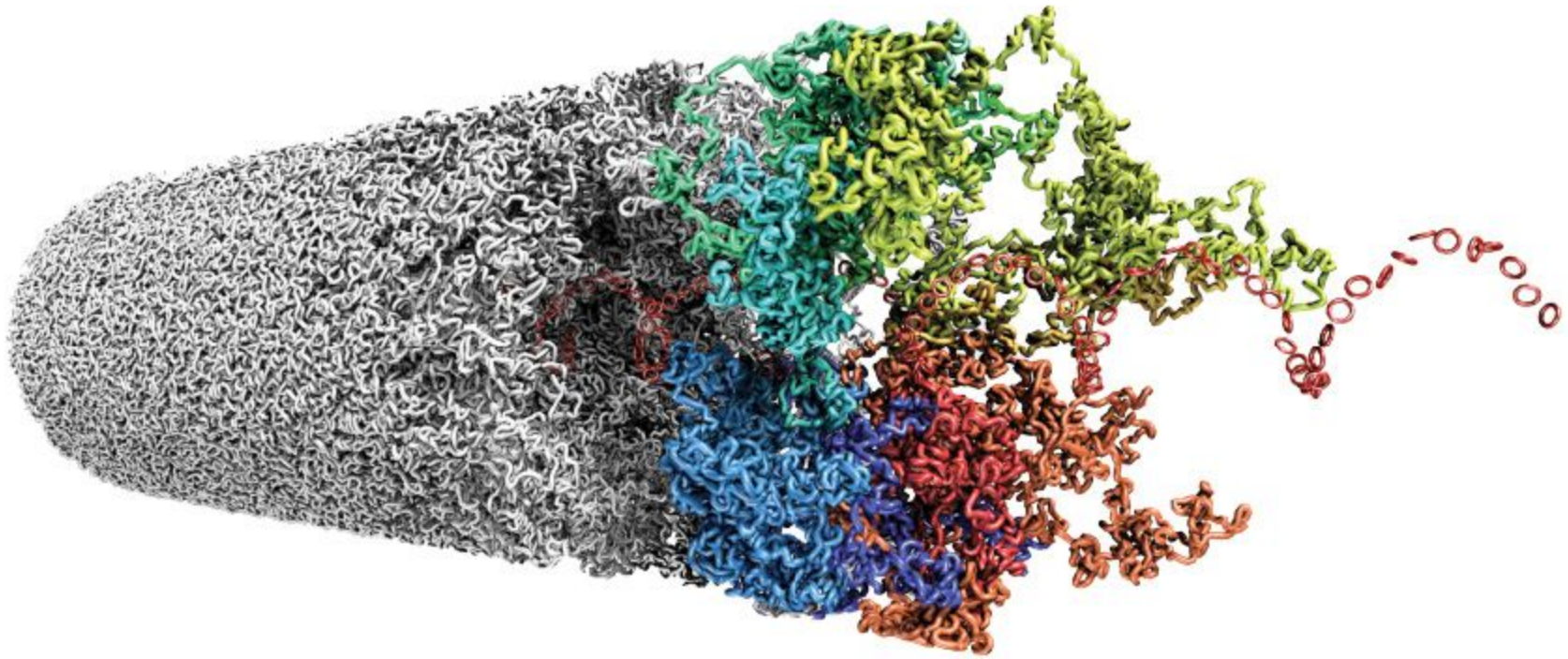
<https://www.sciencedirect.com/topics/neuroscience/rna-seq>

Problems only statistic approaches are facing (ctd.)

- Deficiencies when it comes to robustness and stability across data sets and approaches, due to the statistical issue of working in a high-dimensional space with only a few samples; **Low Robustness and Stability**
- The more complex an approach is, the less trustworthy it becomes, due to users considering them a black box, especially when machine learning techniques are used; **Lack of Transparency**
- Genes sharing similar functions, participating in the same pathway, often share similar expression patterns, \Rightarrow dependencies in expression profiles; **Redundancy**

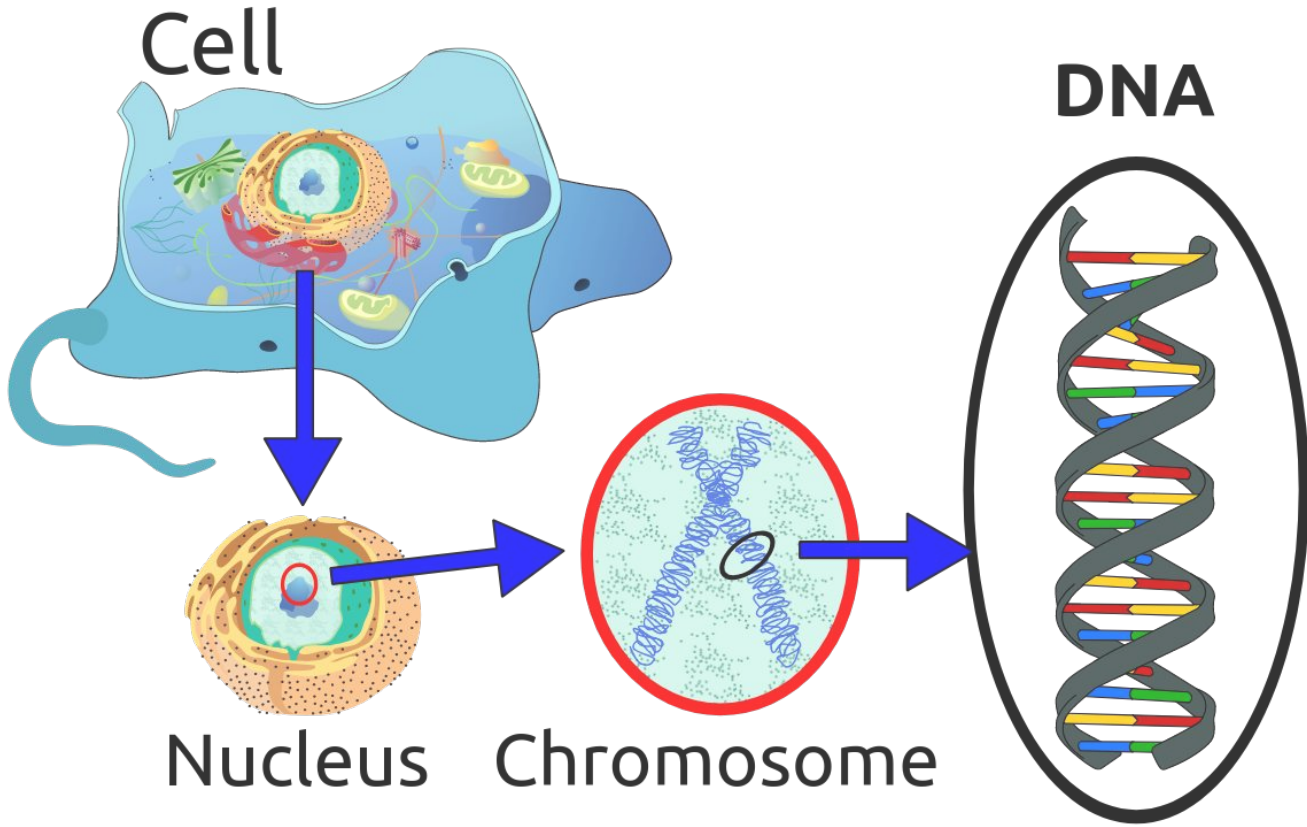


Source: <http://www.vcbio.science.ru.nl/en/virtuallessons/cellcycle/trans/>

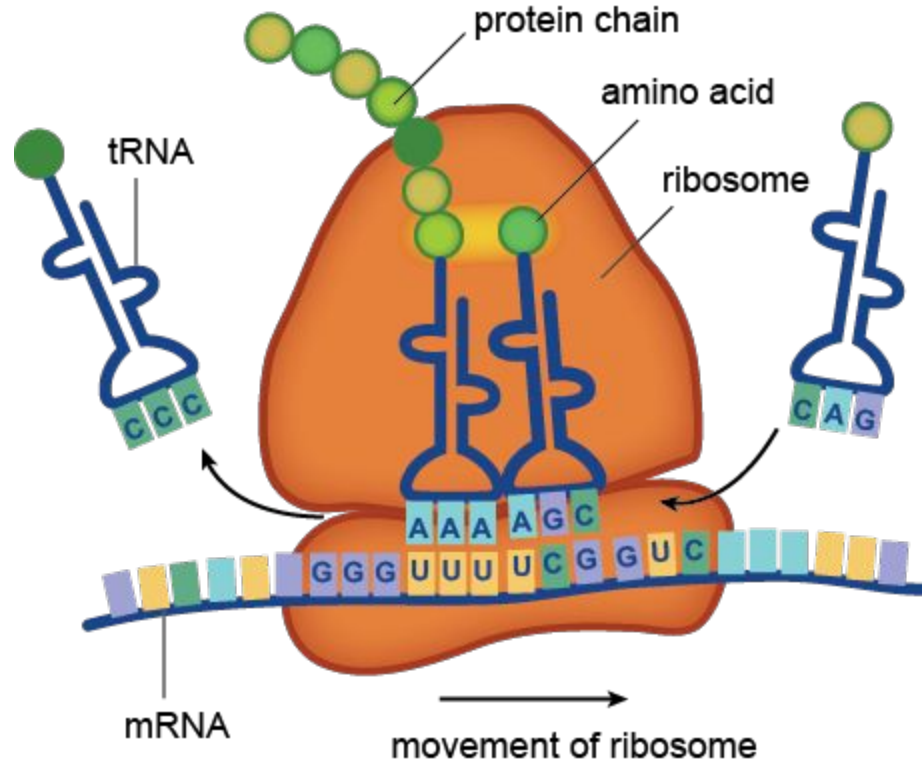


Source.

<https://www.sciencenews.org/article/dna-chromosome-bundles-cell-division-mitosis>



Source: https://da.wikipedia.org/wiki/Fil:Eukaryote_DNA-en.svg



Source: <https://dumielauxepices.net/wallpaper-3773168>

Challenges for integrative gene selection

Genes can be involved in multiple pathways; **Pathway Overlaps**

- Gap between static interaction networks and dynamic cell processes;
Conceptual Gap
- Knowledge bases heavily vary in quality of information, ranking scores, review status of information, evidence, disease specificity and many more aspects;
Knowledge Base Quality

Challenges for integrative gene selection

- It is a fine line between RNAseq data results and integrated external knowledge; the most advanced approach for integrative gene selection has no added value if it cannot be applied in practice due to confusing setup-processes, specific use-cases and so on; updates of knowledge bases sometimes require source-code maintenance / refactorings of integrative gene selection approaches; **Gene Selection Process**
- Current knowledge on biological processes and diseases is not exhaustive, thus knowledge bases are updated frequently, resulting in reproducibility issues as time goes by; **Evaluation**

Motivation for this specific topic / approach

- Untersch. Integrative Ansätze miteinander vergleichen, Network-Based fehlt
- Network-Based Integrative Gene Selection achieved rece
- Anbindung weiterer Knowledge-Bases an das bestehende Framework, um ein generalisierbares Framework für verschiedenste Ziele zu erstellen und Wissenschaftlern an die Hand zu geben, um Ihre eigenen Ansätze mittels dieses Frameworks bewerten zu lassen (Vielleicht einen eigenen Ansatz, der div. Ansätze kombiniert, entwickeln?)

ToDo's:

- Shorten up text @ bulletpoints, only use main keywords
- When a new variant of slides is ready, send to Cindy

Das hier wird wahrscheinlich zu viel...

Genome sizes are typically given as gametic nuclear DNA contents ('C-values') either in units of mass (picograms, where $1 \text{ pg} = 10^{-12} \text{ g}$) or in number of base pairs (in eukaryotes, most often in megabases, where $1 \text{ Mb} = 10^6 \text{ bases}$). These are directly interconvertible as $1 \text{ pg} = 978 \text{ Mb}$ (or $1 \text{ Mb} = 1.022 \times 10^{-3} \text{ pg}$) ([7](#)).

Genome Size of humans: 3.2 Gbp

Genes: approx. 20k-25k

Mitosis & Meiosis - Cell Division (Eukaryotes vs. Gamete)