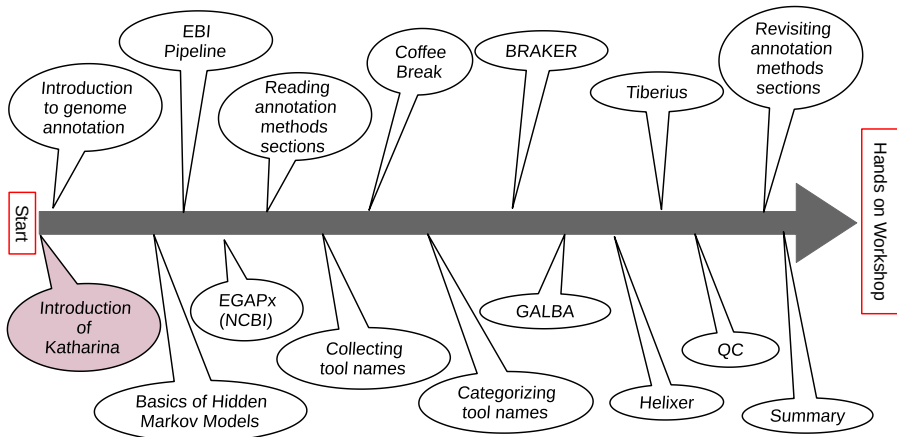# Annotation of Protein Coding Genes

January 8th 2025

Katharina J. Hoff

Twitter: @katharina_hoff
Bluesky: @katharinahoff.bsky.social
Mastodon: @KatharinaHoff@fosstodon.org

E-Mail: katharina.hoff@uni-greifswald.de

# Katharina J. Hoff
Group Leader in Applied Bioinformatics at University of Greifswald

## Short CV

| | |
|---|---|
| 2005 | B.Sc. Plant Biotechnology (Hanover, stays abroad: Budapest & Alnarp) |
| 2009 | Ph.D. Molecular Biology (Göttingen) |
| 2022 | Habilitation (Greifswald) |

## Research

- eukaryotic genome annotation, metagenomics
- best known for: **BRAKER** & other **Gaius-Augustus** software
- 37 peer-reviewed research articles with currently 7,186 citations

## Teaching

- currently 1 postdoc, 4 PhD students, 1 MSc student, 2 BSc students
- applied bioinformatics, programming, statistics, & data science

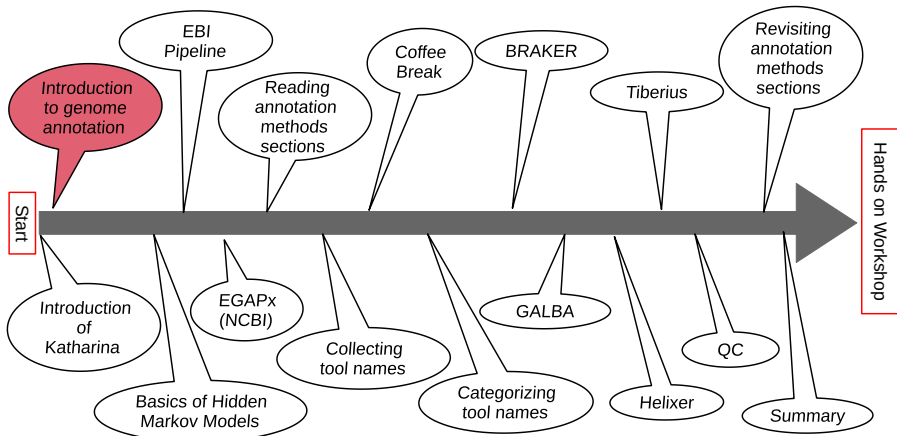... I love to sail, have a dog, a cat, and an 8-years old daughter...

## Materials at

https:
//github.com/KatharinaHoff/GenomeAnnotation_Workshop

(Some images have been removed on Github because I do not have permission to share them.)

## Where are the protein coding genes?
### Genomic sequence: chicken

```
cctcacctctgagaaaacctctttgccaccaataccatgaagctctgcgtgactgtcctgtctctcctc
gtgctagtagctgccttctgctctctagcactctcagcaccaagtaagtctacttttgcagctgctatt
tcgagtcaaggtgtaggcagagtccttttttctagtcatggctggcaaacagtgggatctgggatggg
acaaaaggcagctaggaagattgccatgtagtctgctgctaaatgtagagtctagtagatattcagtaa
cattcaagttcctattttcttaagaattagcaaccagcagaggaaaacgatgggctggaagtcagactg
ttgaattggctctgcctttaattatttgttcaagcaagcccctgtccctctctgtgccttggtttcccc
atctgtcatatgaagggagtgcgatgtgttctgagactgaatccagttccaatcttctagatttctttc
tcgttcttctctgaagatccactattcagaataagactcctgctcatgttaggtgggaatggatacaag
ggaccatatttgggggttctggtagctccacagggatgctcaatgaagatgcaaaattagaagtcaaaat
aaacagctcccatgggcagtgttgatctcaccctggcctttcctttcagtgggctcagaccctcccacc
gcctgctgcttttcttacaccgcgaggaagcttcctcgcaactttgtggtagattactatgagaccagc
agcctctgctcccagccagctgtggtgtgagtatcaacccctggctgccctgggaggcaagggtgaggg
ctggattttttaaaggggggcctgttttggggaggggggtgatgagcgctggggaggcagctctcagggctg
aagccttccctgacagcagtgaggtcacaggtcatgaactcacttttcaagtgctgaaggcggctgagt
ggcagccgagacagaagggggttcctggggaggaagttattcagaggacagggaagcaggggaaggcag
acaggtcccatgagatatggaccaattccttaaaccatgctagaaaaacatgtggaaaagtcactacca
ggctggcagggaatggggcaatctattcatactgattgcaatgcccactggttcctaatctgggcaacc
cctggggcccacagctaaatccagtgagtggaagttacagggagtctgcttccagtgctgctcgaggaa
ggatcccatccaccagagctgccccacatggaccatggtcaggcagaggaagatgcctaccacaggcaa
gggataaagccagatgacctcaaaggtcccatggattctaatctgtctgctccttgttctacagattc
caaaccaaaagaggcaagcaagtctgcgctgaccccagtgagtcctgggtccaggagtacgtgtatgac
ctggaactgaactgagctgctcagagacaggaagtcttc
```

# Examples for the importance of genome annotation

## Silencing polygalacturonase activity in tomato



Sheeny et al. (1988) Proc. Natl. Acad. Sci. USA 85:8805-8809; Image: adapted from

http://luisbarbosa2.blogspot.com/2013/06/flavr-savr-tomato.html, Original: Asia Datta, Subhra Chakraborty, National Institute of Plant

Genome Research, New Delhi

# Examples for the importance of genome annotation
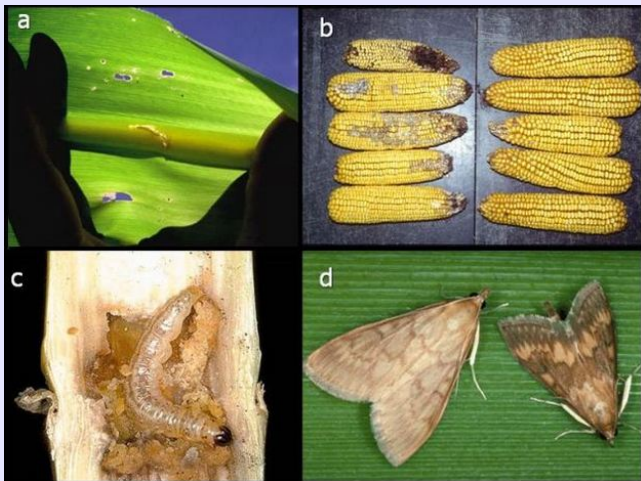
## Bacillus thuringiensis toxin against European corn borer



Image: Hellmich & Hellmich (2012) Nature Education Knowledge 3(10):4

http://www.nature.com/scitable/content/ne0000/ne0000/ne0000/ne0000/46977030/1_2.jpg
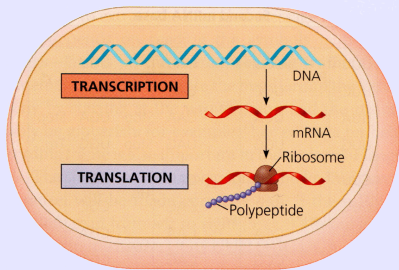
# It does not take a village to publish a genome!

- In the past:
  - ▶ Human: International Human Genome Sequencing Consortium (2001), Nature 409(6822), 860 **248 authors**
  - ▶ Mosquito: Nene et. al (2007) **95 authors**

# It does not take a village to publish a genome!

- In the past:
  - ▶ Human: International Human Genome Sequencing Consortium (2001), Nature 409(6822), 860 **248 authors**
  - ▶ Mosquito: Nene et. al (2007) **95 authors**
- More recently:
  - ▶ 4 *Botrytis cinerea*: Adhikari et al. (2025), **5 authors**
  - ▶ European harvest mouse: O'Brien & Colom (2024), **2 authors**
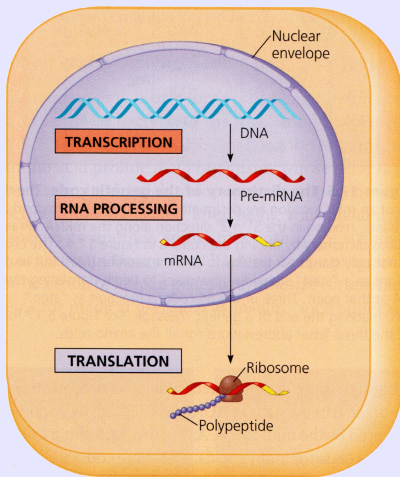  - ▶ Great wood-rush: Goodwin et al. (2024), **4 authors**

# It does not take a village to publish a genome!

- In the past:
  - Human: International Human Genome Sequencing Consortium (2001), Nature 409(6822), 860 **248 authors**
  - Mosquito: Nene et. al (2007) **95 authors**
- More recently:
  - 4 *Botrytis cinerea*: Adhikari et al. (2025), **5 authors**
  - European harvest mouse: O'Brien & Colom (2024), **2 authors**
  - Great wood-rush: Goodwin et al. (2024), **4 authors**
- **You can do it!**

# How does a cell recognize protein-coding genes?
## Transcription & Translation



Images: Campbell et al., Biology, San Francisco, 2008, p. 329, Fig. 17.3

# How does a cell recognize protein-coding genes?
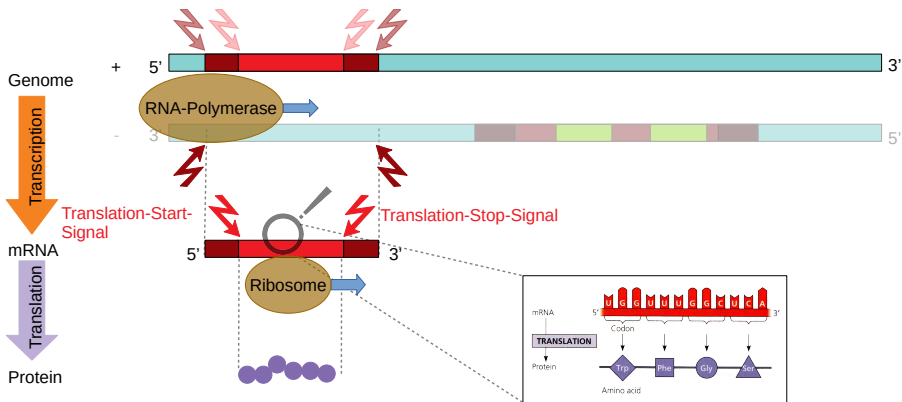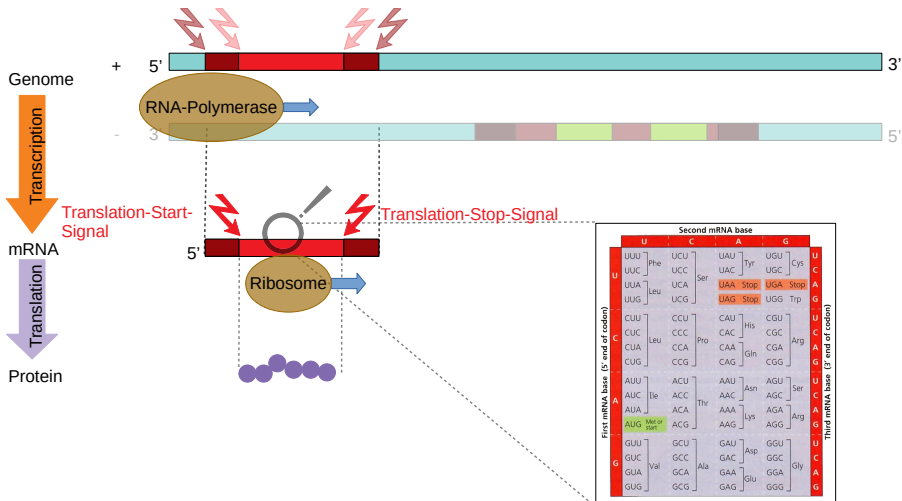## Prokaryotes & Eukaryotes*



Image: Campbell et al., Biology, San Francisco, 2008, p. 329, Fig. 17.4

*) only some of the genes in eukaryotes

# How does a cell recognize protein-coding genes?
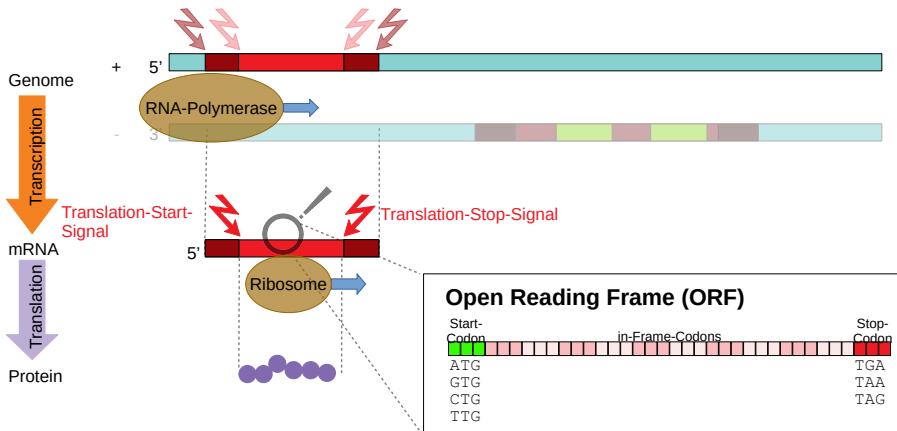## Prokaryotes & Eukaryotes*



Image: Campbell et al., Biology, San Francisco, 2008, p. 339, Fig. 17.5

*) only some of the genes in eukaryotes
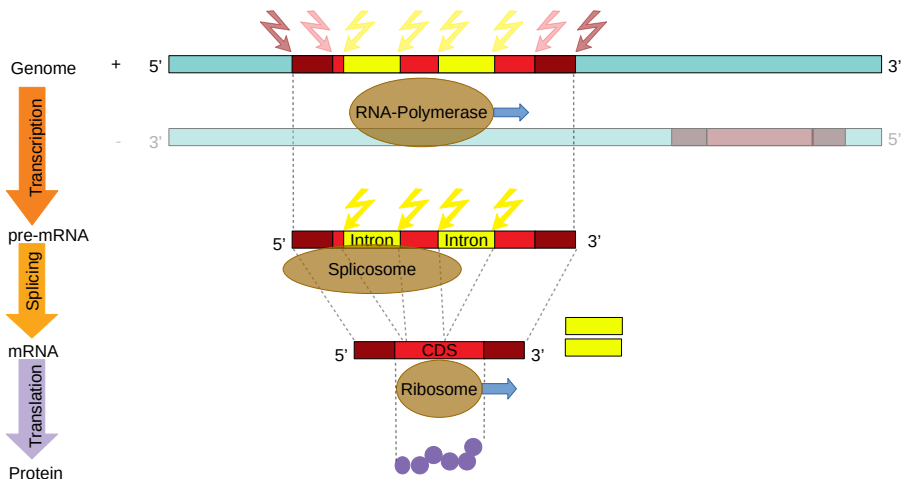
# How does a cell recognize protein-coding genes?
## Prokaryotes & Eukaryotes*



- every protein coding gene has an ORF
- not every ORF is a protein coding gene

# How does a cell recognize protein-coding genes?
Eukaryotes: Splicing of introns

# The Genome Annotation Problem
## Genomic Sequence: chicken

```
cctcacctctgagaaaacctctttgccaccaataccatgaagctctgcgtgactgtcctgtctctcctc
gtgctagtagctgccttctgctctctagcactctcagcaccaagtaagtctacttttgcagctgctatt
tcgagtcaaggtgtaggcagagtccttttttctagtcatggctggcaaacagtgggatctggggatggg
acaaaaggcagctaggaagattgccatgtagtctgctgctaaatgtagagtctagtagatattcagtaa
cattcaagttcctattttcttaagaattagcaaccagcagaggaaaacgatgggctggaagtcagactg
ttgaattggctctgcctttaattatttgttcaagcaagcccctgtccctctctgtgccttggtttcccc
atctgtcatatgaagggagtgcgatgtgttctgagactgaatccagttccaatcttctagatttctttc
tcgttcttctctgaagatccactattcagaataagactcctgctcatgttaggtgggaatggatacaag
ggaccatatttgggttctggtagctccacagggatgctcaatgaagatgcaaaattagaagtcaaaat
aaacagctcccatgggcagtgttgatctcaccctggcctttcctttcagtgggctcagaccctcccacc
gcctgctgcttttcttacaccgcgaggaagcttcctcgcaactttgtggtagattactatgagaccagc
agcctctgctcccagccagctgtggtgtgagtatcaacccctggctgccctgggaggcaagggtgaggg
ctggattttaaaggggggcctgttttggggaggggtgatgagcgctggggaggcagctctcagggctg
aagccttccctgacagcagtgaggtcacaggtcatgaactcacttttcaagtgctgaaggcggctgagt
ggcagccgagacagaagggggttcctggggaggaagttattcagaggacagggaagcaggggaaggcag
acaggtcccatgagatatggaccaattccttaaaccatgctagaaaaacatgtggaaaagtcactacca
ggctggcagggaatggggcaatctattcatactgattgcaatgcccactggttcctaatctgggcaacc
cctggggcccacagctaaatccagtgagtggaagttacagggagtctgcttccagtgctgctcgaggaa
ggatcccatccaccagagctgccccacatggaccatggtcaggcagaggaagatgcctaccacaggcaa
gggataaagccagatgacctcaaaggtcccatggattctaatctgtctgctccttgttctacagattc
caaaccaaaagaggcaagcaagtctgcgctgaccccagtgagtcctgggtccaggagtacgtgtatgac
ctggaactgaactgagctgctcagagacaggaagtcttc
```
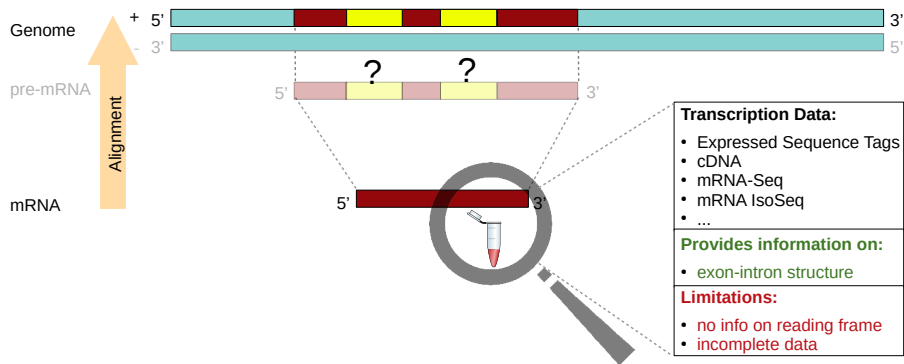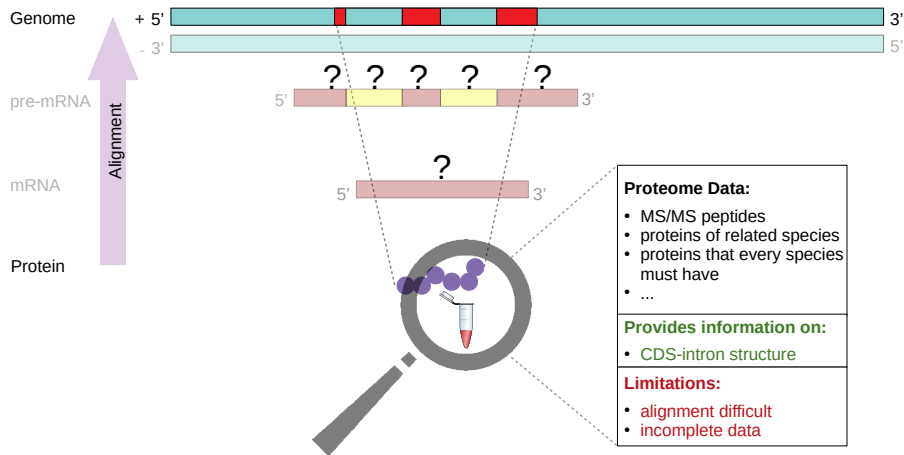
# The Genome Annotation Problem

Genomic sequence: chicken (1 gene: macrophage inflammatory protein-1 b)

```
cctcacctctgagaaaacctctttgccaccaataccatgaagctctgcgtgactgtcctgtctctcctc
gtgctagtagctgccttctgctctctagcactctcagcaccaagtaagtctacttttgcagctgctatt
tcgagtcaaggtgtaggcagagtccttttttctagtcatggctggcaaacagtgggatctggggatggg
acaaaaggcagctaggaagattgccatgtagtctgctgctaaatgtagagtctagtagatattcagtaa
cattcaagttcctattttcttaagaattagcaaccagcagaggaaaacgatgggctggaagtcagactg
ttgaattggctctgcctttaattatttgttcaagcaagcccctgtccctctctgtgccttggtttcccc
atctgtcatatgaagggagtgcgatgtgttctgagactgaatccagttccaatcttctagatttctttc
tcgttcttctctgaagatccactattcagaataagactcctgctcatgttaggtgggaatggatacaag
ggaccatatttggggttctggtagctccacagggatgctcaatgaagatgcaaaattagaagtcaaaat
aaacagctcccatgggcagtgttgatctcaccctggcctttcctttcagtgggctcagaccctcccacc
gcctgctgcttttcttacaccgcgaggaagcttcctcgcaactttgtggtagattactatgagaccagc
agcctctgctcccagccagctgtggtgtgagtatcaacccctggctgcctgggaggcaagggtgaggg
ctggatttttaaaggggggcctgttttggggagggggtgatgagcgctggggaggcagctctcagggctg
aagccttccctgacagcagtgaggtcacaggtcatgaactcacttttcaagtgctgaaggcggctgagt
ggcagccgagacagaaggggttcctggggaggaagttattcagaggacagggaagcaggggaaggcag
acaggtcccatgagatatggaccaattccttaaaccatgctagaaaaacatgtggaaaagtcactacca
ggctggcagggaatggggcaatctattcatactgattgcaatgcccactggttcctaatctgggcaacc
cctgggggcccacagctaaatccagtgagtggaagttacagggagtctgcttccagtgctgctcgaggaa
ggatcccatccaccagagctgccccacatggaccatggtcaggcagaggaagatgcctaccacaggcaa
gggataaagccagatgacctcaaaggtcccatgggattctaatctgtctgctccttgttctacagattc
caaaccaaaagaggcaagcaagtctgcgctgacccagtgagtcctgggtccaggagtacgtgtatgac
ctggaactgaactgagctgctcagagacaggaagtcttc
```

# What aids in the identification of genes in genomes?
## Evidence data from transcription

# What aids in the identification of genes in genomes?
## Evidence data from translation

# What aids in the identification of genes in genomes?
## Mathematical models



Genome

**Mathematical models:**
- **Hidden Markov Models**
  (e.g. GeneMark, AUGUSTUS)
- dynamic programming
- Support Vector Machines
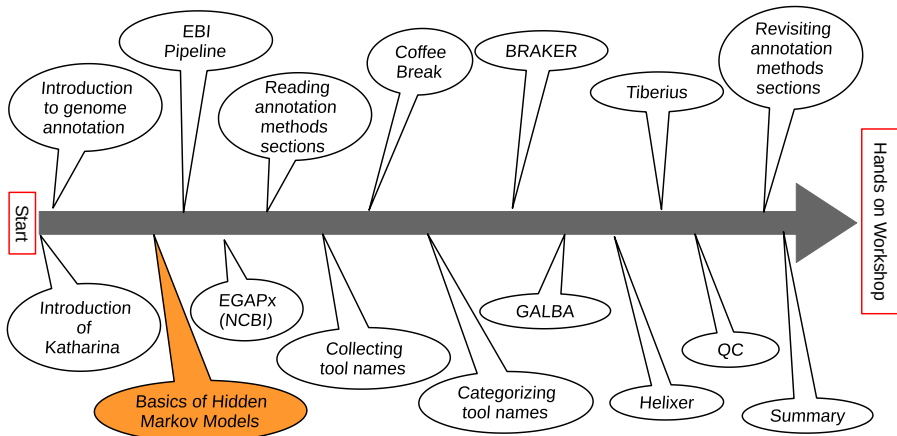- neural networks
- decision tree systems
- ...

**Provide information on:**
- complete gene structures (sometimes incl. UTRs)

**Limitations**
- *predictions* may be wrong
- models use **parameters** that have to be trained

**How does the cell recognize?**
- transcription signals
- translation signals
- splicing signals
- ORF

# What aids in the identification of genes in genomes?
## Mathematical models



Genome

**Mathematical models:**
- **Hidden Markov Models** (e.g. GeneMark, AUGUSTUS)
- dynamic programming
- Support Vector Machines
- neural networks
- decision tree systems
- ...

**Provide information on:**
- complete gene structures (sometimes incl. UTRs)

**Limitations**
- *predictions* may be wrong
- models use **parameters** that have to be trained

A **Hidden Markov Model** can read the genome sequence from left to right and, through knowledge of signals for transcription and translation, assign a probable state to each nucleotide (e.g., intergenic region or CDS).

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

- There are only 2 nucleotides: A, B
- There are only 2 sequence states: intergenic (I), coding sequence (K)

**Input: "Genome sequence"**

e.g. AABBBAB

**Goal: "Most likely path through hidden states"**

```
e.g.  AABBBAA
or    IIKKIKI     P(path) = 0.3%
```

# Basis of highly accurate gene prediction tools
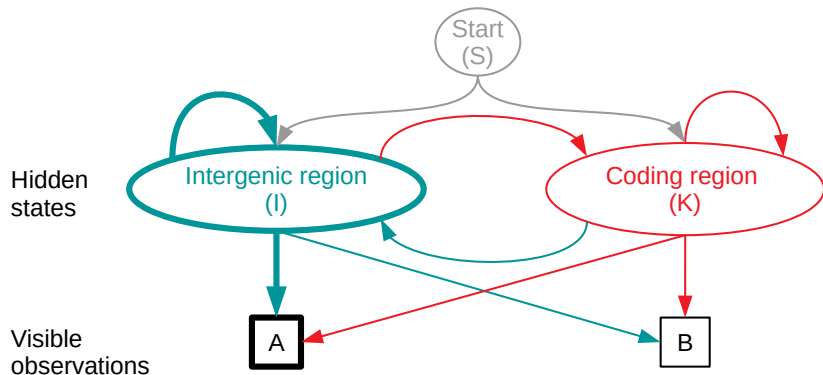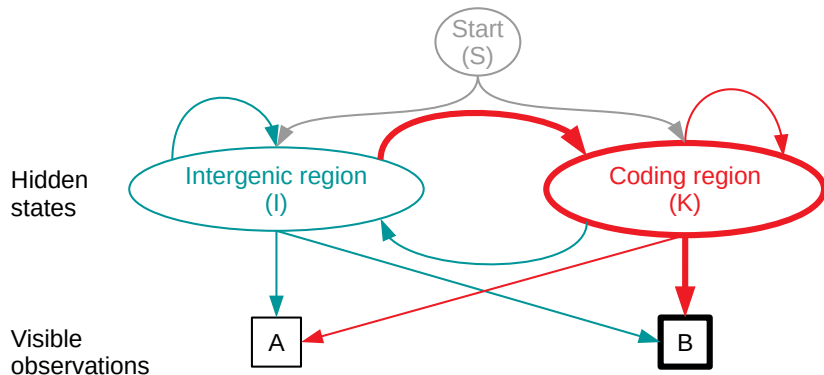## Hidden Markov Model



Hidden states

**A possible 'state path' for the genome sequence:**

AABBBAA

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

**A possible 'state path' for the genome sequence:**

AABBBAA
I

# Basis of highly accurate gene prediction tools
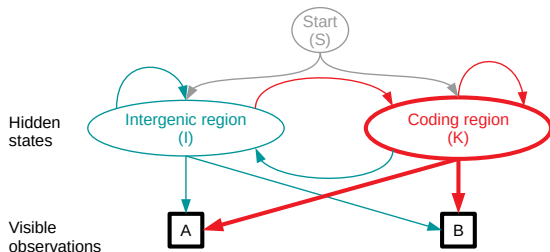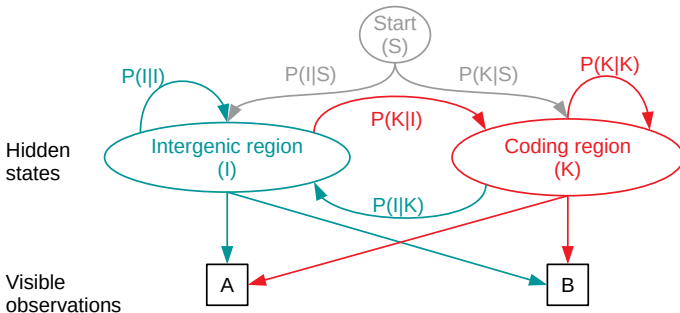## Hidden Markov Model



Hidden states

**A possible 'state path' for the genome sequence:**

AABBBAA
II

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**A possible 'state path' for the genome sequence:**

AABBBAA
IIK...

## Model properties

1. The current value of the hidden state depends exclusively on the state of its predecessor.

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**A possible 'state path' for the genome sequence:**

**A possible 'state path' for the genome sequence:**

A
I

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**A possible 'state path' for the genome sequence:**

AA
II

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**A possible 'state path' for the genome sequence:**

```
AAB...
IIK...
```

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



## Model properties

1. The current value of the hidden state depends exclusively on the state of its predecessor.
2. The current value of the visible observation depends exclusively on the value of the current, hidden state.

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**How likely are the state transitions?**

# Basis of highly accurate gene prediction tools
Hidden Markov Model

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

# Basis of highly accurate gene prediction tools
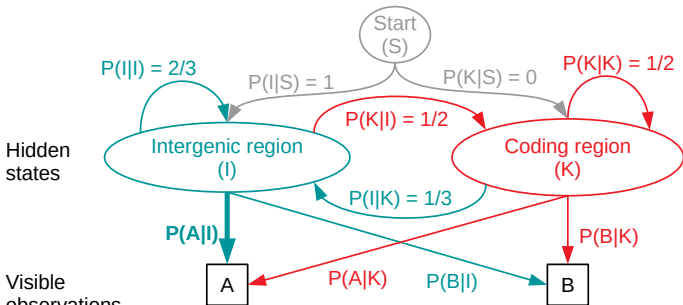## Hidden Markov Model



**Training data:**

AABABA
IKKIII
   -++

P(I|I) = 2/3          P(I|K) = 1 – P(I|I) = 1/3

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

Start (S)

$P(I|I) = 2/3$   $P(I|S) = 1$   $P(K|S) = 0$   $P(K|K) = 1/2$

$P(K|I) = 1/2$

Hidden states

Intergenic region (I)   Coding region (K)

$P(I|K) = 1/3$

**P(A|I) = 3/4**   $P(B|K)$

Visible observations (emissions)

A   $P(A|K)$   **P(B|I) = 1/4**   B

**How likely are the observations?**

```
AABABA
IKKIII
+  +-+
```

```
P(A|I) = ¾          P(B|I) = 1 – P(A|I) = 1 – ¾ = ¼
```

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**How likely is a given state-emission path?**

```
Path = AAB
       IKK

P(Path) = ?
```

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



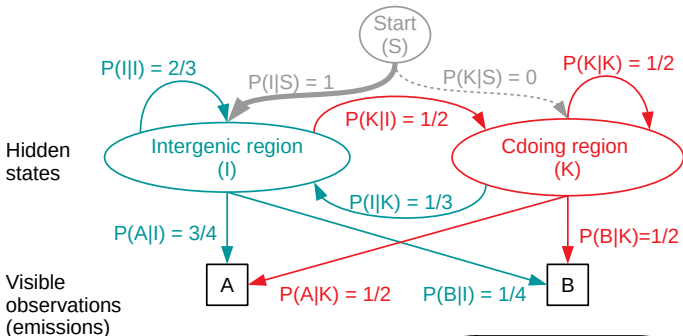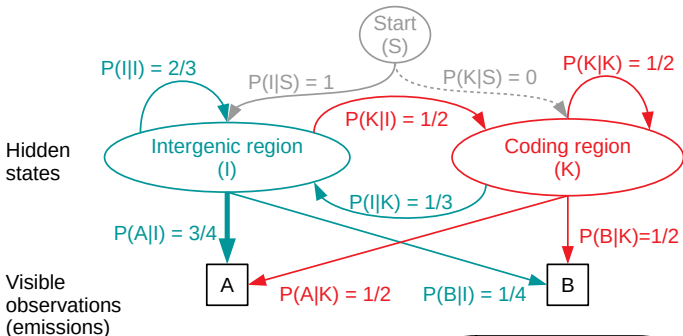**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

P(Path) = P(I|S)

Multiply the probabilities along the state-emission path!

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



Start
(S)

P(I|I) = 2/3

P(I|S) = 1

P(K|S) = 0

P(K|K) = 1/2

P(K|I) = 1/2

Hidden
states

Intergenic region
(I)

Coding region
(K)

P(I|K) = 1/3

P(A|I) = 3/4

P(B|K)=1/2

Visible
observations
(emissions)

A

P(A|K) = 1/2

P(B|I) = 1/4

B

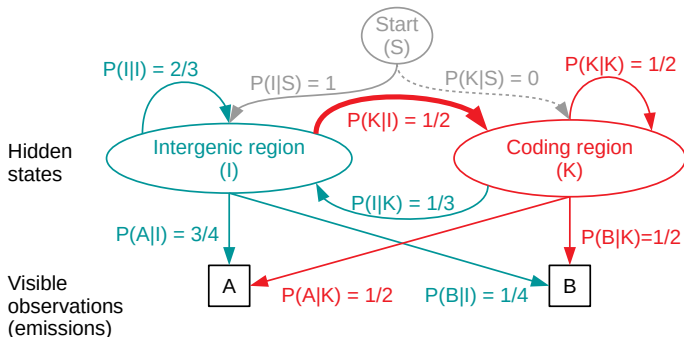**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

$P(Path) = P(I|S)*P(A|I)$

Multiply the probabilities
along the state-emission
path!

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



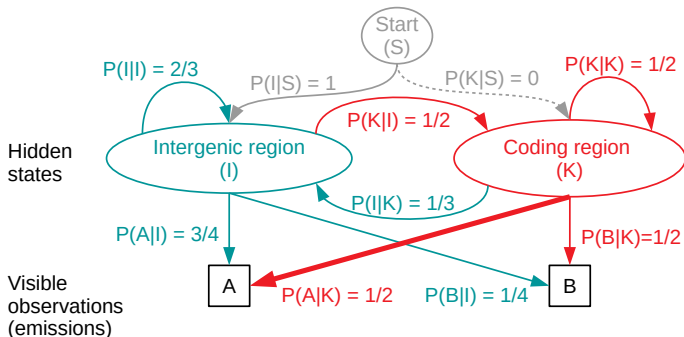**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

$$P(Path) = P(I|S)*P(A|I)*P(K|I)$$

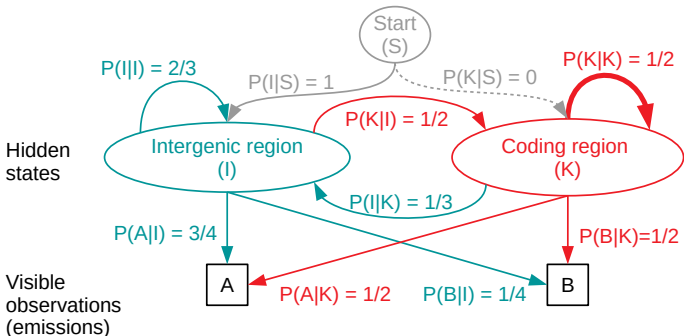# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

$P(Path) = P(I|S)*P(A|I)*P(K|I)*P(A|K)$

# Basis of highly accurate gene prediction tools
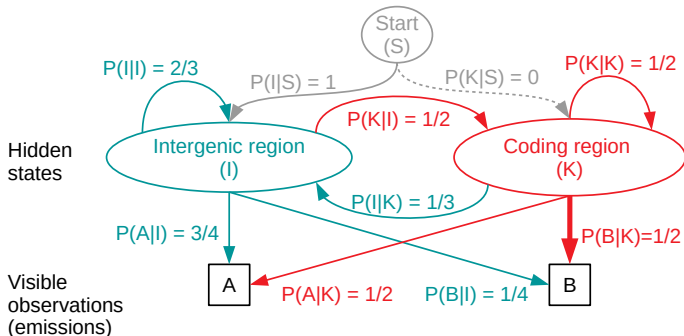## Hidden Markov Model



**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

$P(Path) = P(I|S)*P(A|I)*P(K|I)*P(A|K)*P(K|K)$

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



**How likely is a given state-emission path?**

```
Path = AAB
       IKK
```

P(Path) = P(I|S)*P(A|I)*P(K|I)*P(A|K)*P(K|K)*P(B|K)

# Basis of highly accurate gene prediction tools
## Hidden Markov Model



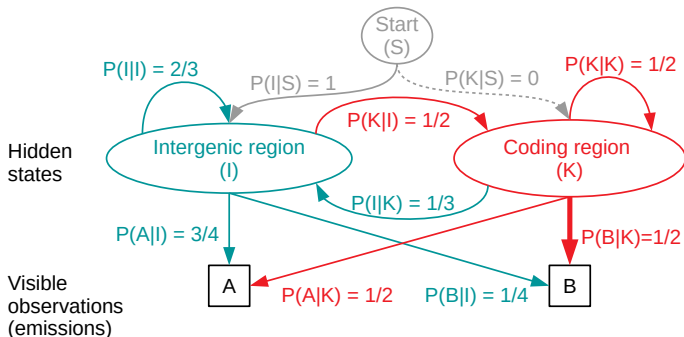**How likely is a given state-emission path?**

```
Path = AAB
       IKK

P(Path) = P(I|S)*P(A|I)*P(K|I)*P(A|K)*P(K|K)*P(B|K)

        = 1 * ¾ * ½ * ½ * ½ * ½

        = 3/64
```

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

## Find the most probable state sequence for a given sequence

**Input: "genome sequence"**

AABBBABA

**Problem: "too many possible state sequences"**

IIIKKKKKK
KKIKKIIIK
IIKIIIKIK
IKKIKIIIK
KIKIKKKIK
KKKIKIKKK
...

Idea:

1. Generate all possible state sequences
2. Calculate the probability for each state sequence
3. Choose the state sequence with the highest probability

$\Rightarrow$ too expensive!

# Basis of highly accurate gene prediction tools
## Hidden Markov Model

Find the most probable state sequence for a sequence: Viterbi Algorithm.
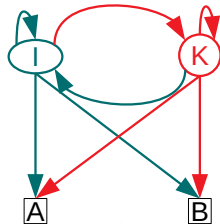
Transition probabilities

Emission probabilities
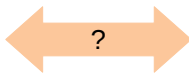
AABBBABA

Viterbi

Most probable state sequence:
IIKKKIIII

- 4096 observed nucleotide hexamers
- Many more hidden states
  (e.g. 3'-UTR, 5'-UTR, Intron, ...)

Gene
(Parameter training)

?

Gene
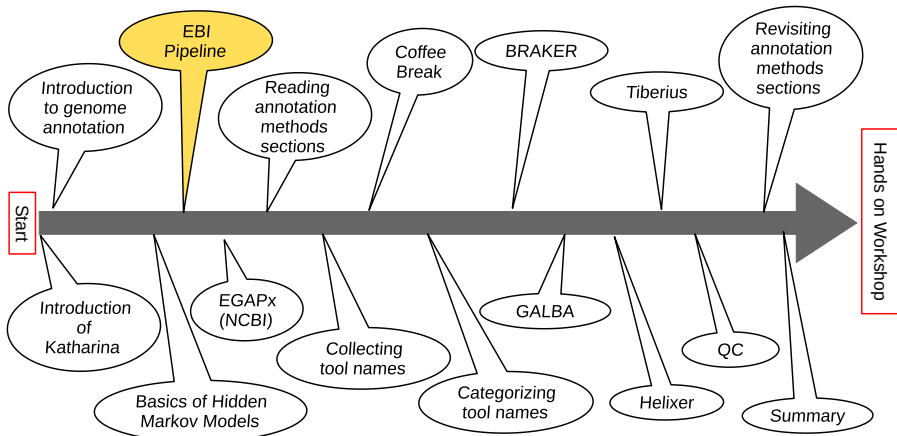(Prediction)

Genome

Mathematical models:
- **Hidden Markov Models**
  (e.g. GeneMark, AUGUSTUS)
- dynamic programming
- Support Vector Machines
- neural networks
- decision tree systems
- ...

**Provide information on:**
- complete gene structures (sometimes incl. UTRs)

**Limitations**
- *predictions* may be wrong
- models use **parameters** that have to be trained

**Transcription data:**
- Expressed Sequence Tags
- cDNA
- mRNA-Seq
- mRNA IsoSeq
- ...

**Proteome data:**
- MS/MS peptides
- proteins of related species
- proteins that every species must have
- ...

# EBI: Ensembl annotation system

# Ensembl annotation pipelines

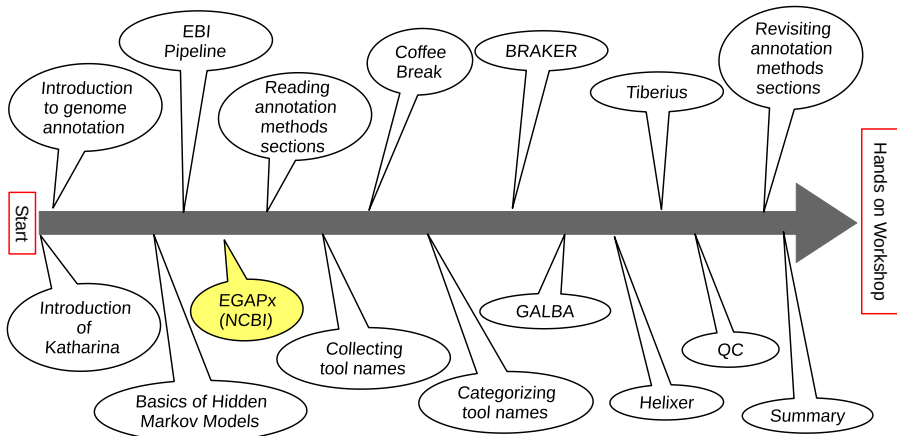# Ensembl annotation pipeline for non-vertebrates

# Genome Annotation at the NCBI

- Internal: The NCBI Eukaryotic Genome Annotation Pipeline (EGAP)
- (Internal: RefSeq curation)
- You can run it: **EGAPx**

# Annotation with EGAPx (NCBI)

- Containerized with Docker/Singularity
- Documentation: `https://github.com/ncbi/egapx`
- Currently supported clades (protein sets):
  - ▶ Chordata
  - ▶ Insecta
  - ▶ Arthropoda
  - ▶ Monocots
  - ▶ Eudicots
- Easy to use
- Benchmarking possible: good accuracy!

# Read your methods snippet
Focus on structural annotation of protein coding genes only!

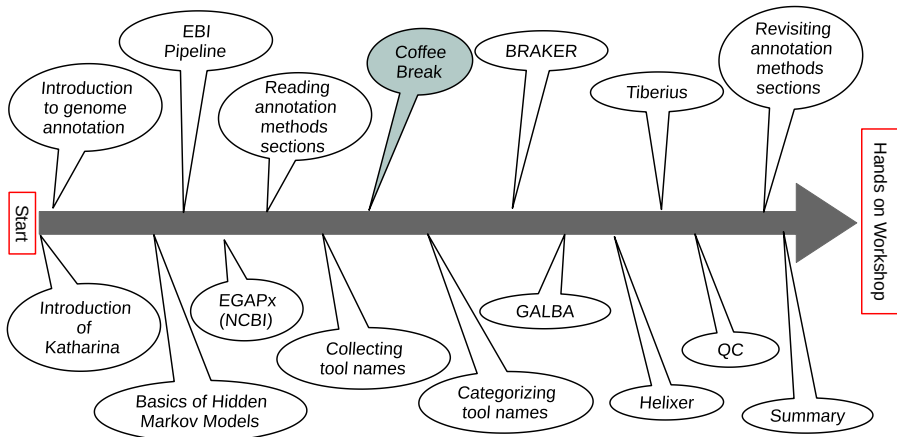1. We move to Wooclap
2. Enter the names of tools involved in structural annotation of protein coding genes
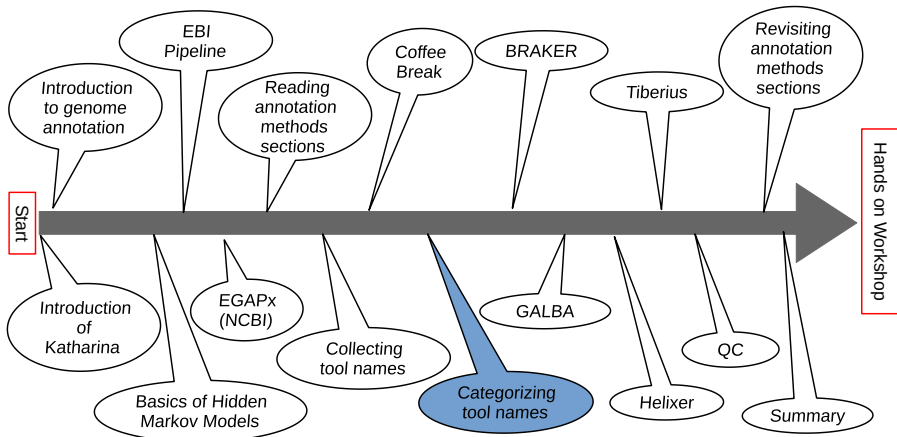
# Read your methods snippet
## Focus on structural annotation of protein coding genes only!

Names of tools involved in structural annotation of protein coding genes

# Categorize tool names

Go to

`https://shorturl.at/uA0Tg`

and sort the tools names from your methods snippet into categories

# Categorize tool names

## Protein-to-Genome Aligners +

**Diamond**
Not a spliced aligner, a fast mapper. Does not directly return evidence for gene finding but can be used to select candidates for accurate spliced alignment to save runtime. Also used for removing redundancy in training gene sets (in protein-protein mode).

**TBLASTN**
Not a spliced aligner, a fast mapper. Does not directly return evidence for gene finding. Not much used in the field anymore because DIAMOND is faster.

**GenomeThreader**
Not used much anymore because miniprot is faster and more robust towards increasing distance between donor and recepient

**miniprot**
very fast

**Spaln**
Spaln3 is also very fast

**(GeMoMa)**
Does not require protein sequence as input but the genomes and annotations of closely related species. Returns accurate gene models if mapping is successful.

**Exonerate**
very slow

**(EviAnn)**
see note on GeMoMa, not related tool but works in a similar way in terms of possible inputs.

## Transcriptome Processing Tools +

**HISAT2**
Fast and efficient short reads spliced aligner.

**STAR**
Remember to run 2-pass mapping.

**Minimap**
For long reads

**StringTie**
State of the art and fast transcriptome assembler (genome-guided assembly). Can use short and/or long reads.

**Trinity**
For de novo transcriptome assembly. This is often not very helpful for structural genome annotation and needs a lot of runtime...

**PASA**
Old tool but still pretty good, includes Transdecoder for finding ORFs in assembled transcripts, can perform UTR annotation

**cDNA cupcake**
Outdated.

**cd-hit**
Used to cluster transcripts, not directly helpful for genome annotation but sometimes during data preparation.

## ab initio Gene Finders +

**AUGUSTUS**
Older but still very accurate gene finder. Can run ab initio or with evidence. Was in the past also used in MAKER, is at the core of BRAKER and Galba.

**GeneMark**
Suite of self-training gene finding tools. GeneMarkS-T: finding genes in transcripts. GeneMarkES: ab initio, genome sequence input only

**Helixer**
deep learning gene finder, game changer in the field, good BUSCO accuracy, poor gene structure accuracy, great web service, several clade models available

**Tiberius**
deep learning gene finder, building on results of Helixer team. So far only trained for mammals (poor parameters for diatoms also exist)

**FgeneSH**
Older gene finder that can run ab initio or with evidence. Was often used in MAKER in the past.

**GlimmerHMM**
Older gene finder, was an early community project. Was in the past often used in MAKER.

**SNAP**
Older gene finder that can run ab initio or with evidence. Was often used in MAKER in the past. Very easy to train and use.

## Gene Finders that Use Evidence +

**PASA**
Can be used to generate training genes for gene finders from transcriptome evidence. Not much used anymore for this.

**GeMoMa**
Only evidence based gene models. Requires as input genomes and annotations of related species. Can add transcriptome data.

**EviAnn**
Similar to GeMoMa

**AUGUSTUS**
Can use lots of evidence. There is also a comparative version (AUGUSTUS-CGP) that annotates genes consistently in a multi-species genome alignment.

**GeneMark**
GeneMarkET: genome +rnaseq mapping input
GeneMarkEP: genome + protein db input
GeneMarkETP: genome+rnaseq mapping+protein db input

**SNAP**

**FGenesH**

## Gene Set Combiners +

**TSEBRA**
Custom tailored to AUGUSTUS/GeneMark output and evidence, strong focus on splice site support

**EVM**
Widely used, kind of maximizes coverage information for building an optimal gene set

**(MAKER)**
contains a combiner part to build a consensus gene set

# Categorize tool names

## Gene Set Combiners

**TSEBRA**
Custom tailored to AUGUSTUS/GeneMark output and evidence, strong focus on splice site support

**EVM**
Widely used, kind of maximizes coverage information for building an optimal gene set

**(MAKER)**
contains a combiner part to build a consensus gene set

## Complex Annotation Pipelines

**BRAKER**
BRAKER3: RNA-Seq + protein database
BRAKER2: protein database
(BRAKER1: RNA-Seq spliced mapping file, better to use BRAKER3)

**EASEL**
Nextflow pipeline, includes several transcriptome assembles, protein mappers, and gene finders.

**TOGA**
Interesting if you have multi-genome alignment, integrates AUGUSTUS-CGP

**MAKER2**
The first highly popular community annotation pipeline. Now outdated because accuracy is not that good. Does not train gene finders automatically.

**PGAP (prokaryotes)**

**Prokka (prokaryotes)**

## Functional Genome Annotation Tools

**InterProScan**

**Blastp**

**EggNOG-mapper**

**CD-SEARCH**

**InterProScan**

**FANTASIA**
Very fast, only GO term assignment

## Repeat Masking Tools

**RepeatMasker**

**RepeatModeler2**

**HEletronScanner**

**Deep TE**

**LTR FINDER**

**GMATA**

**RepeatScout**

**Tandem Repeat Finder**

**EDTA pipeline**

**Red**

## Others

**BUSCO**
Marker gene detection in geomes, proteomes and transcriptomes. Widely used.

**OMArk**
Marker gene detection in proteomes. Handles alternative splicing isoforms well. Larger number of marker genes than in BUSCO.

**AGAT**
Useful for many gff handling tasks

**OrthoDB**
Database

**Gffcompare**
Useful for handling gff files

Image: credits to DALL-E2, modified by human

## The BRAKER Team
University of Greifswald & Georgia Tech University



Lars Gabriel

Alexandre Lomsadze, Katharina Hoff, Tomáš Brůna
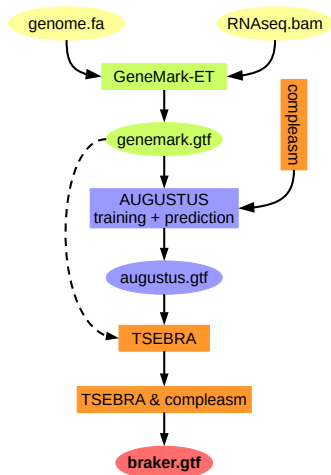
Mario Stanke

Mark Borodovsky

Also: Simone Lange, Matthis Ebel, Hannah Thierfeldt, Anica Hoppe, Neng Huang

# BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS 🆓

Katharina J. Hoff ✉, Simone Lange, Alexandre Lomsadze, Mark Borodovsky ✉, Mario Stanke

- spliced alignments of RNA-Seq are used by GeneMark-ET and AUGUSTUS
- 1,677 citations (Google Scholar)

### Whole-Genome Annotation with BRAKER

**Katharina J. Hoff, Alexandre Lomsadze, Mark Borodovsky, and Mario Stanke**

in Kollmar M. (eds) Gene Prediction. Methods in Molecular Biology, vol 1962. Humana, New York, NY, 2019

# GeneMark-ET uses RNA-Seq for Training

## Anchors from RNA-Seq for training



**Figure 3.** Selection of elements of training set in GeneMark-ET for the next iteration. The new training set of protein-coding regions is comprised from exons with at least one 'anchored splice site' as well as long exons predicted *ab initio* (>800 nt).

- employs unsupervised training
- training includes introns and exons anchored by mapped RNA-Seq reads
- does not require RNA-Seq reads assembly
- does not use RNA-Seq information in the *prediction* step

AUGUSTUS uses RNA-Seq for **Prediction**

Introns predicted by RNA-Seq read alignment

Genome

AUGUSTUS gene predictions with "hints" from RNA-Seq

- requires "prior data" for training
- uses intron information from RNA-seq for *prediction*
- no RNA-Seq assembly required
- optional input: BUSCO lineage (compleasm)

# Measuring accuracy of genome annotation

## Experiments

Accuracy assessment after applying tool to genome with reference annotation:

| Species | Genome Size (Mb) | # Genes in Annotation |
|---|---|---|
| *Arabidopsis thaliana* (thale cress) | 119 | 27,444 |
| *Bombus terrestris* (bumble bee) | 249 | 10,581 |
| *Caenorhabditis elegans* (nematode) | 100 | 20,172 |
| *Danio rerio* (zebrafish) | 1,345 | 25,611 |
| *Drosophila melanogaster* (fruit fly) | 137 | 13,928 |
| *Gallus gallus* (chicken) | 1,040 | 17,279 |
| *Medicago truncatula* (barrelclover) | 420 | 44,464 |
| *Mus musculus* (mouse) | 2,650 | 22,378 |
| *Parasteatoda tepidariorum* (house spider) | 1,445 | 18,602 |
| *Populus trichocarpa* (poppy) | 389 | 34,488 |
| *Solanum lycopersicum* (tomato) | 772 | 33,562 |

## Accuracy metrics

**Precision** = Specificity: Percentage of correctly found genes/transcripts/exons in the **predicted gene set**.

**Recall** = Sensitivity: Percentage of correctly found genes/transcripts/exons in the **reference annotation**.

**F1-Score**: $\dfrac{2 \cdot \text{Recall} \cdot \text{Precision}}{\text{Recall} + \text{Precision}}$

Use only if not enough RNA-Seq for BRAKER3!

# BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database
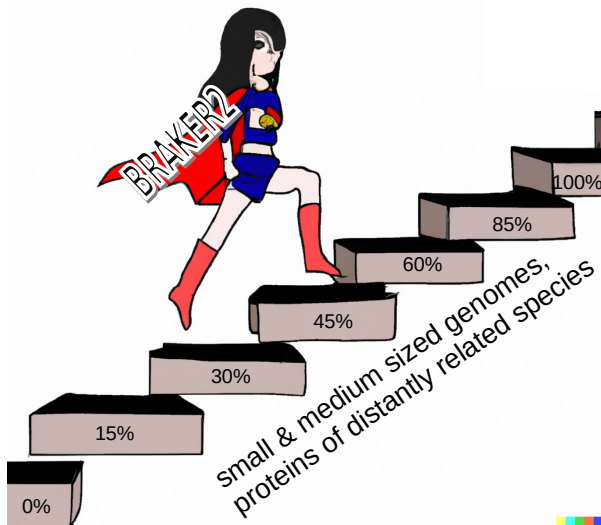
**Tomáš Brůna[1,†], Katharina J. Hoff[2,3,†], Alexandre Lomsadze[4], Mario Stanke[2,3,‡] and Mark Borodovsky** [4,5,*,‡]



- spliced alignments of a large number of proteins (e.g. OrthoDB partition)
- optional input: BUSCO lineage (compleasm)
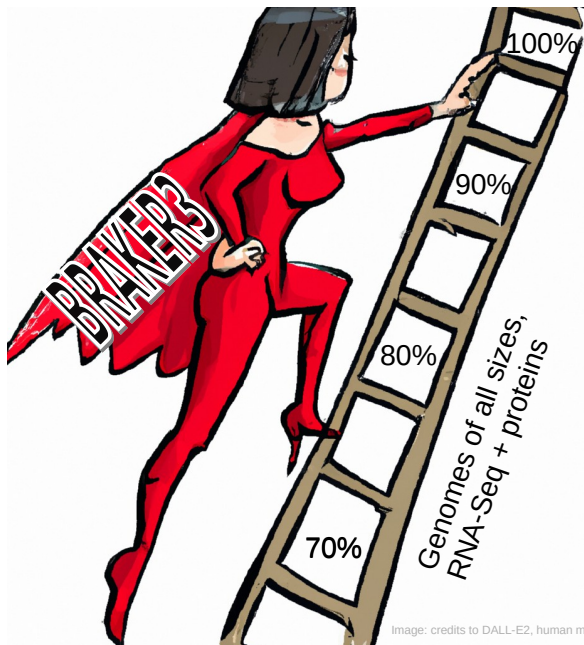- 1,269 citations (Google Scholar)

Use only if you have no RNA-Seq data on genomes <1 Gbp

Image: credits to DALL-E2, human modification

# BRAKER3: using RNA-Seq and protein evidence with GeneMark-ETP, AUGUSTUS and TSEBRA



- Gabriel *et al.* (2024)
- 137 citations (Google Scholar)
- spliced aligned and **assembled** RNA-Seq
- large protein database
- optional input: BUSCO lineage (compleasm)
- combines GeneMark-ETP and AUGUSTUS gene sets with TSEBRA

# BRAKER3: using RNA-Seq and protein evidence with GeneMark-ETP, AUGUSTUS and TSEBRA



HISAT2: Kim et al. (2019). Nature Biotechnology, 37(8), 907-915.
StringTie: Pertea et al. (2015). Nature Biotechnology, 33(3), 290-295.
GeneMarkS-T: Tang et al. (2015). Nucleic Acids Research, 43(12), e78-e78.
GeneMark-ETP: Brůna et al. (2024). Genome Research.
AUGUSTUS: Mario Stanke, et al., 2008. Bioinformatics. 24.5:637-644.
Compleasm: Huang & Li (2023). Bioinformatics, 39(10), btad595.
TSEBRA: Lars Gabriel, et al., 2021. BMC Bioinformatics. 22:566.

**SOFTWARE** **Open Access**

# TSEBRA: transcript selector for BRAKER

Lars Gabriel[1,2], Katharina J. Hoff[1,2], Tomáš Brůna[3], Mark Borodovsky[4,5] and Mario Stanke[1,2*]

- **combines** several gene sets according to evidence
- 154 citations (Google Scholar)



Can be used to combine BRAKER1 and BRAKER2 output if BRAKER3 fails.

**Figure 2.** Average precision and sensitivity of gene predictions made by BRAKER1, BRAKER2, TSEBRA, GeneMark-ETP, and BRAKER3 for the genomes of 11 different species (listed in Supplemental Table S1). Inputs were the genomic sequences, short-read RNA-seq libraries, and protein databases (*order excluded*).

Image: Gabriel *et al.* (2024), Genome Research

# Availability

## Docker/Singularity

```
singularity build braker.sif \
        docker://teambraker/braker:latest

singularity exec braker.sif braker.pl [OPTIONS]
```

## Licenses

- BRAKER: Artistic License
- most components unter open source software licenses
- GeneMark-ETP: CC BY-NC

# GALBA Contributors



Tomáš Brůna   Heng Li   Joseph Guhlin   Lars Gabriel   Natalia Nenasheva

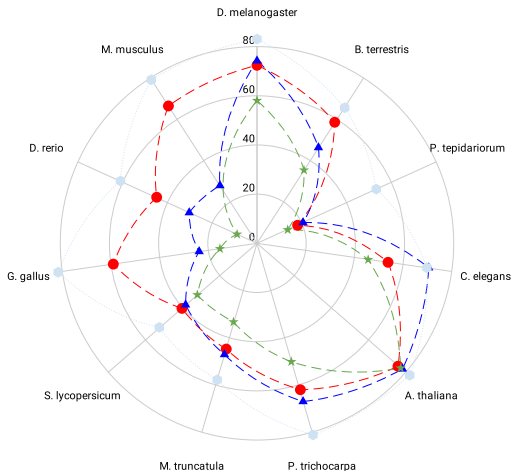Ethan Tolman   Paul Frandsen   Matthis Ebel   Mario Stanke   Katharina Hoff

## Galba: genome annotation with miniprot and AUGUSTUS

- 28 citations (Google Scholar)
- for genomes >1Gbp
- proteins of close relatives

# Proteins Only (GALBA, BRAKER2, FunAnnotate) vs. BRAKER3 with RNA-Seq & Proteins



Gene F1 (%)

**If you have RNA-Seq, use BRAKER3!**

# Availability

## GitHub

```
https://github.com/Gaius-Augustus/GALBA
```

## Docker/Singularity

```
singularity build galba.sif \
        docker://katharinahoff/galba:latest

singularity exec galba.sif galba.pl [OPTIONS]
```

## Licenses

- GALBA: Artistic License
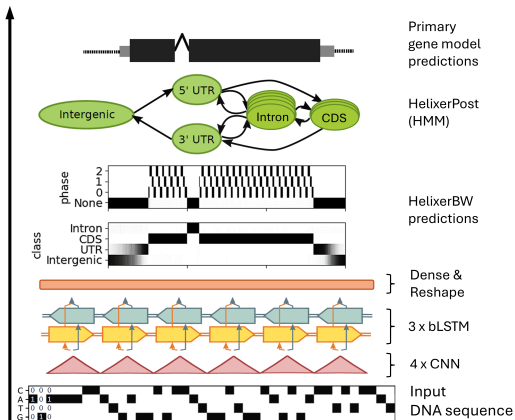- all dependencies have Open Source Licenses

# Helixer: bringing deep learning into genome annotation



Image: ChatGPT by OpenAI, manual editing

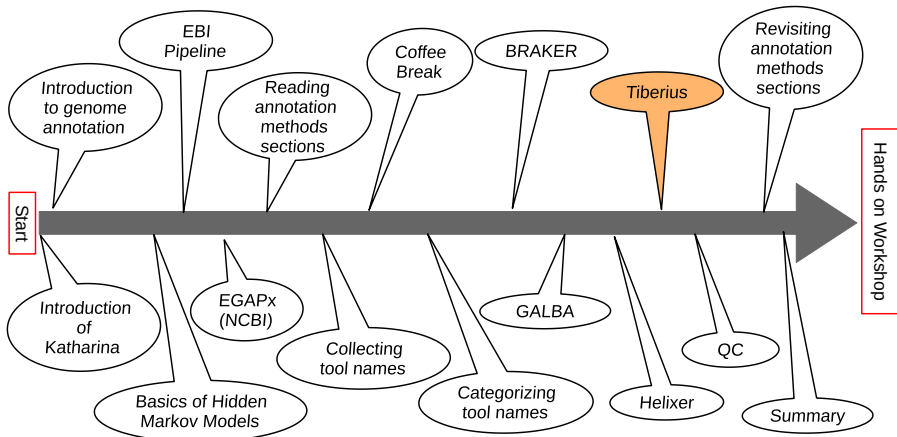# HELIXER–*de novo* PREDICTION OF PRIMARY EUKARYOTIC GENE MODELS COMBINING DEEP LEARNING AND A HIDDEN MARKOV MODEL

Felix Holst[1†], Anthony Bolger[2†], Christopher Günther[1], Janina Maß[3],
Sebastian Triesch[1,4], Felicitas Kindel[1], Niklas Kiel[1,4], Nima Saadat[3,4], Oliver Ebenhöh[3,4],
Björn Usadel[2,4,5], Rainer Schwacke[2], Marie Bolger[2], Andreas P.M. Weber[1,4], Alisandra K. Denton[1,4]

- 27 citations (Google Scholar)
- cross-species gene finder
- *ab initio* prediction
- Pre-trained models for:
  - fungi
  - land plant
  - vertebrate
  - invertebrate
- accuracy (BUSCO): good
- web service

**Availability:** https://github.com/weberlab-hhu/Helixer
Image of Helixer: https://github.com/weberlab-hhu/Helixer/blob/main/img/network.png

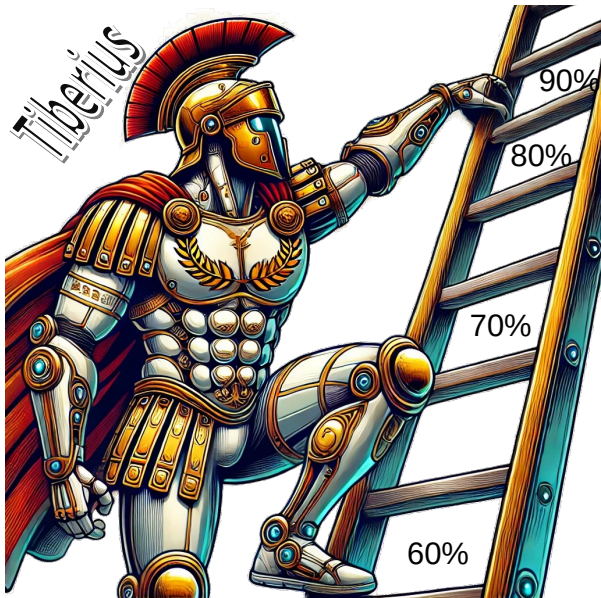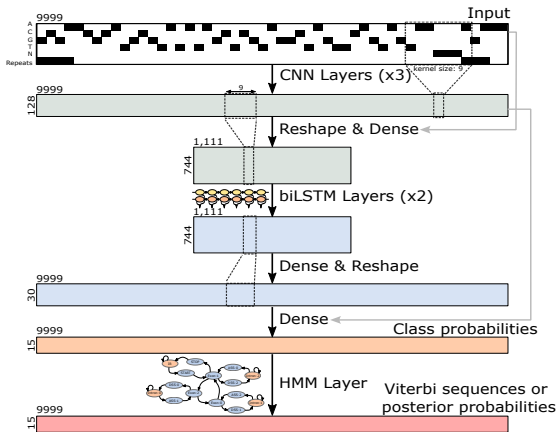# Tiberius: improved genome annotation with deep learning



Image: ChatGPT by OpenAI, manual editing

Lars Gabriel, Felix Becker, Katharina Hoff, Mario Stanke

# Tiberius: end-to-end deep learning with an HMM for gene prediction

**Lars Gabriel** [1,2],∗, **Felix Becker** [1,2], **Katharina J. Hoff** [1,2], **Mario Stanke** [1,2],∗



- builds on findings by Helixer team
- cross-species gene finder
- faster
- higher accuracy
- *ab initio* prediction
- Pre-trained model(s) for:
  - ▶ mammals
  - ▶ (diatoms)
- container for A100 GPU

Availability: https://github.com/Gaius-Augustus/Tiberius

# Accuracy of state of the art gene finders
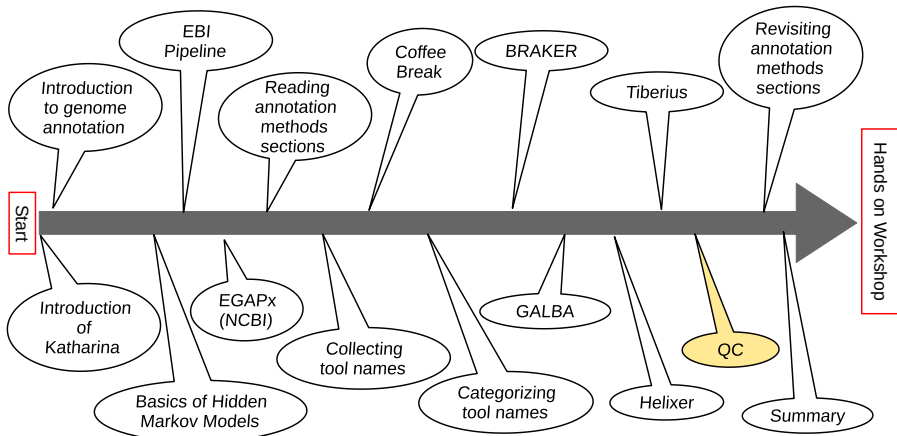## No alternative splicing isoforms

Image: Lars Gabriel, PAG presentation 2025

# Current shortcomings of deep learning gene finders

- no evidence integration
- no alternative splicing isoform prediction
- require expensive GPU for feasible runtime
- limited to specific clades
$\rightarrow$ BRAKER3, Galba & EGAPx currently remain important

# Did We Do a Good Job?

# Genome Browsers
Visualize your Annotation in Context with Evidence

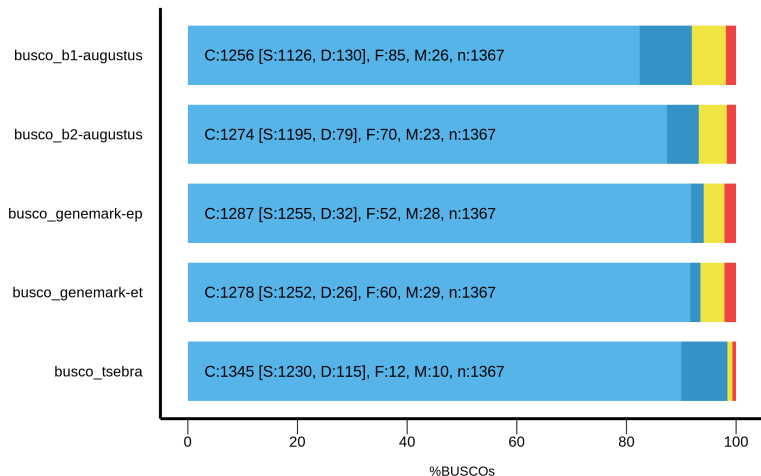- UCSC Genome Browser, MakeHub
- JBrowse
- ...

# Describe Your Annotation

- number of genes
- number of transcripts
- ratio of mono-exonic to multi-exonic genes
- median number of exons per transcript
- maximal number of exons per transcript
- median transcript length
- ...

If possible, compare to annotated close relatives.
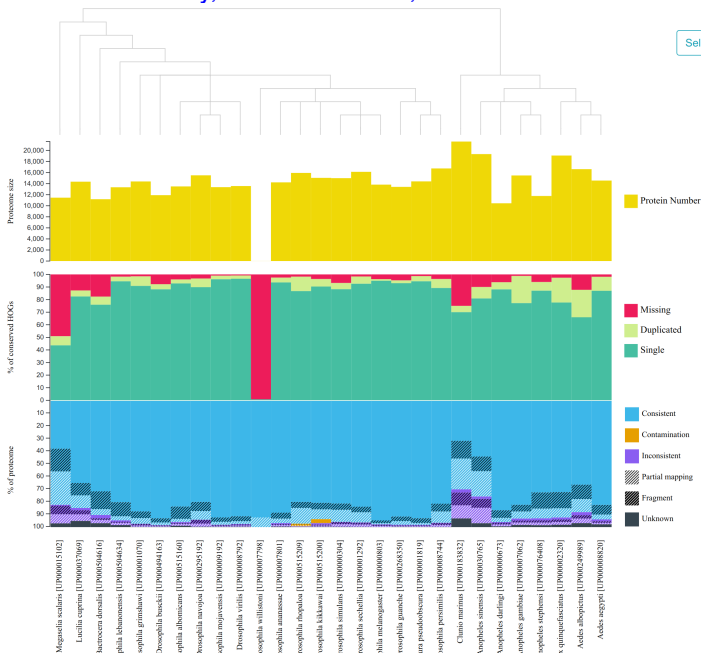Consider effect of individual annotation pipelines.

# BUSCO: Sensitivity in Clade-Specific Conserved Genes



**BUSCO Assessment Results**

- Complete (C) and single-copy (S)
- Complete (C) and duplicated (D)
- Fragmented (F)
- Missing (M)

busco_b1-augustus: C:1256 [S:1126, D:130], F:85, M:26, n:1367

busco_b2-augustus: C:1274 [S:1195, D:79], F:70, M:23, n:1367

busco_genemark-ep: C:1287 [S:1255, D:32], F:52, M:28, n:1367

busco_genemark-et: C:1278 [S:1252, D:26], F:60, M:29, n:1367

busco_tsebra: C:1345 [S:1230, D:115], F:12, M:10, n:1367

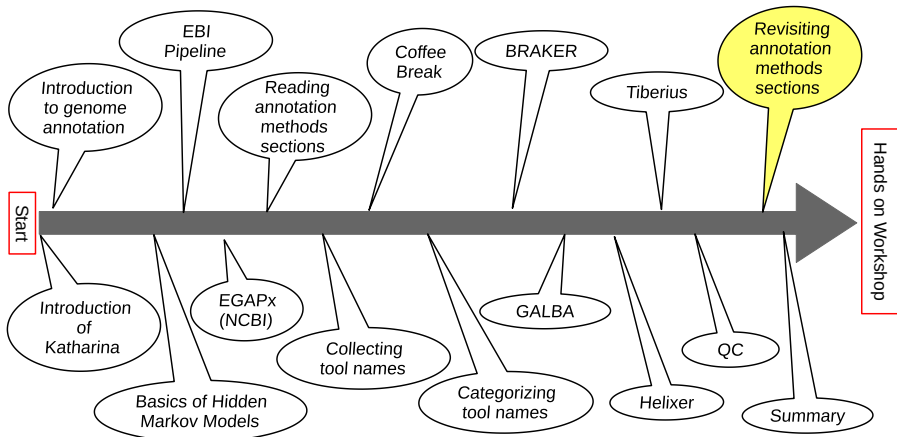%BUSCOs

**Beware!** BUSCO completeness does not warrant correct gene structures!
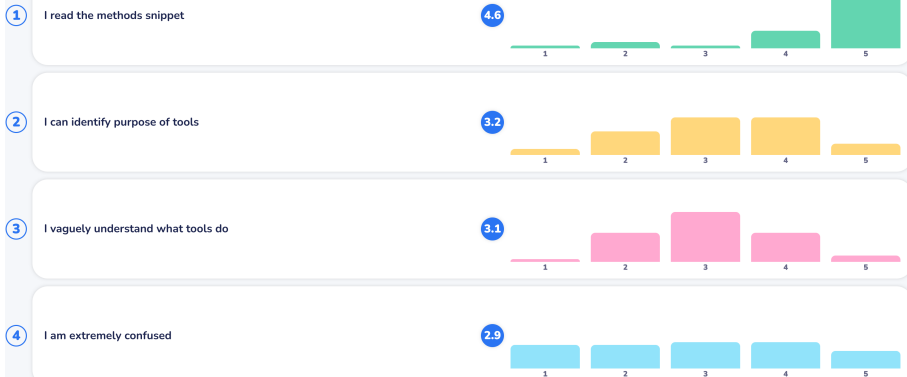
# OMArk: Sensitivity, Contaminations, & More

# Revisiting annotation methods sections
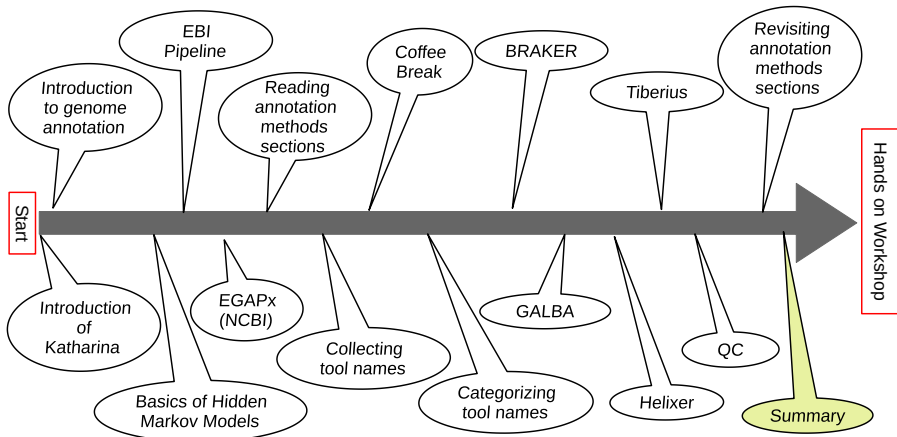
## Your tasks

1. Read your methods snippet, again
2. Use our categorized tool name board at
   `https://shorturl.at/uA0Tg` if you are still unsure what a tool does
3. Ask if you remain unsure what a method is good for
4. Fill the poll on Wooclap

# Revisiting annotation methods sections



Revisiting Annotation Methods (1 = disagree strongly, 5 = completely agree)

① I read the methods snippet — 4.6

② I can identify purpose of tools — 3.2

③ I vaguely understand what tools do — 3.1

④ I am extremely confused — 2.9

# Most important stuff on genome annotation

- structural genome annotation in eukaryotes is hard
- Hidden Markov Models are essential
- evidence helps a lot
- majority of genomes is annotated by large centers
- popular community annotation pipelines:
  1. BRAKER
  2. GALBA
  3. (EGAPx may become popular)
- deep learning is changing the field
  1. Helixer (careful with accuracy)
  2. Tiberius (only for two clades)
- "looking nice" is not always "correct"
- BUSCO completeness is widely used
- OMArk might be more appropriate
- high marker gene detection rate $\neq$ high accuracy