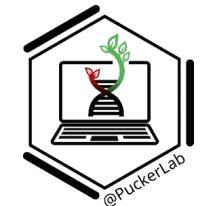
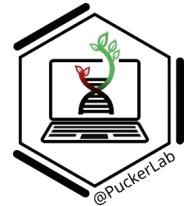
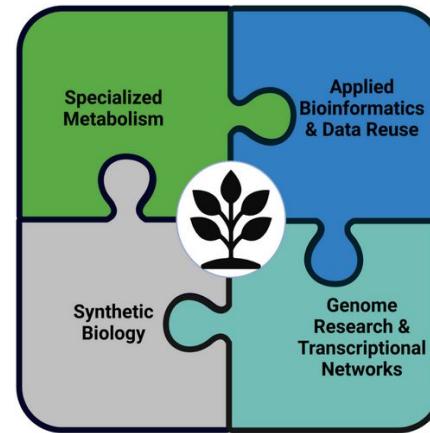


# de.NBI course 2: From Gene Models to Biological Insights - Functional Annotation

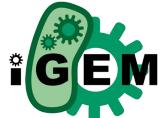
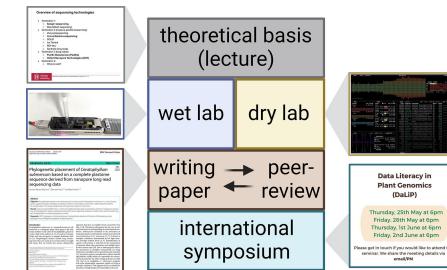




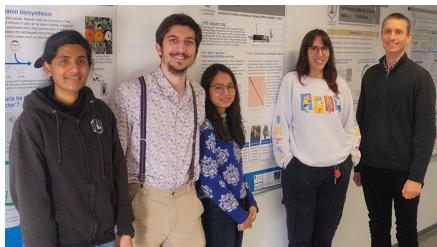
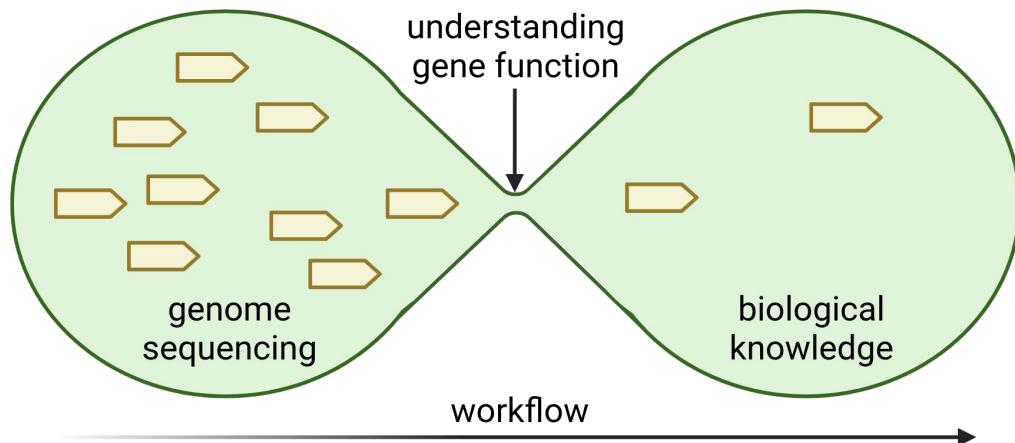
- Flavonoids: Anthocyanins, Flavonols, Flavones, Proanthocyanidins
- Withanolides
- Carotenoids
- ...



Tools: KIPES, NAVIP, MGSE, MYB/bHLH\_annotator, DupyliCate...

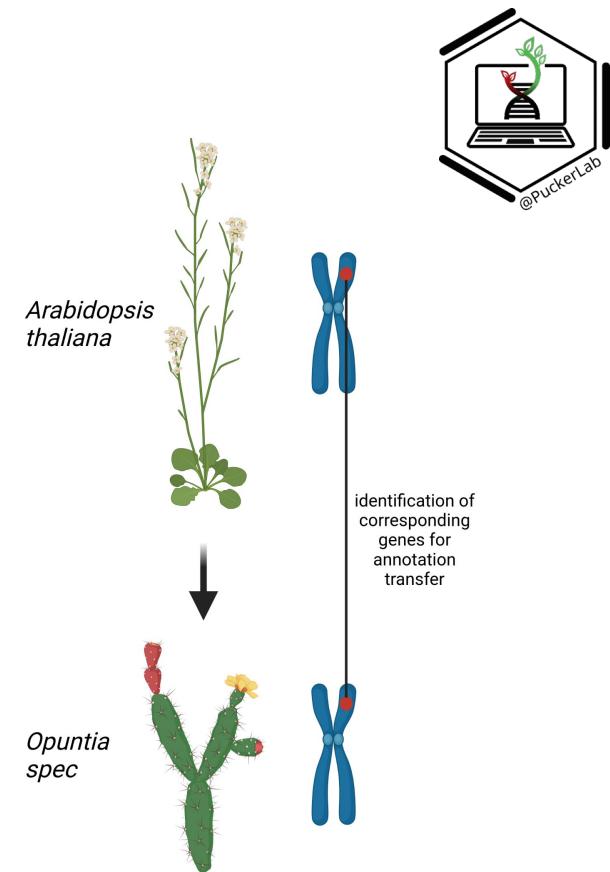


# Course motivation

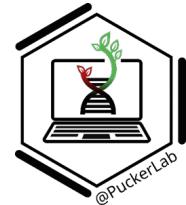


From L to R: Shakunthala, Samuel, Nancy, Julie, and Boas  
Picture credits: Jakob Horz

From Gene Models to Biological Insights - Functional Annotation | Day 1 | 3



## Availability of Slides

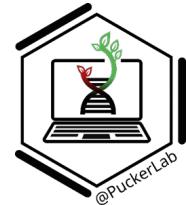


All materials will be shared with participants at the end of this course. Background information is available here:

<https://www.preprints.org/manuscript/202508.1176>

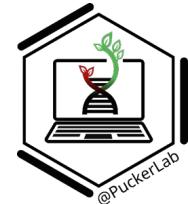
**Slides** → <https://github.com/bpucker/teaching/tree/master/deNBI2025>



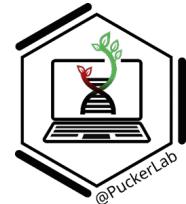


Today you will learn:

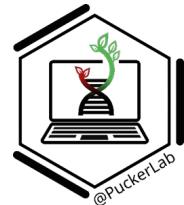
- Introduction to Virtual Machine
- Access to Virtual Machine
- Running BUSCO
- Transposable element annotation



- high performance computing => server farm



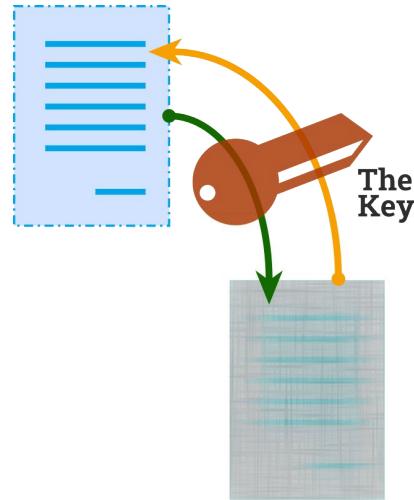
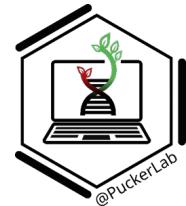
- high performance computing => server farm
- de.NBI offers VMs (VM = virtual machine)
  - running Linux
  - access through SSH
- de.NBI connected LifeScience AAI Account needed



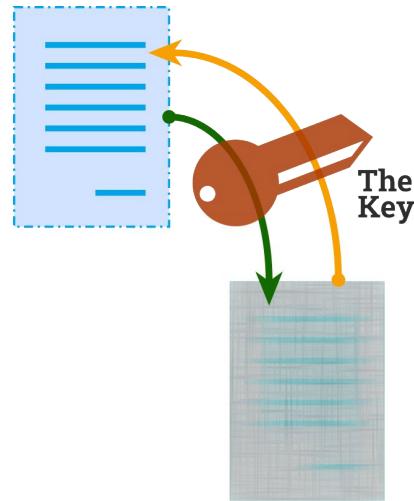
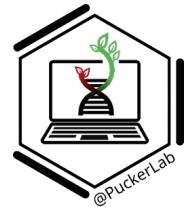
- high performance computing => server farm
- de.NBI offers VMs (VM = virtual machine)
  - running Linux
  - access through SSH
- de.NBI connected LifeScience AAI Account needed
- Tasks:
  1. Register: <https://cloud.denbi.de/wiki/registration/>
  2. Apply for the course de.NBI project:  
<https://signup.aai.lifescience-ri.eu/fed/registrar/?vo=denbi&group=pbbcourse20251>

Who is new?

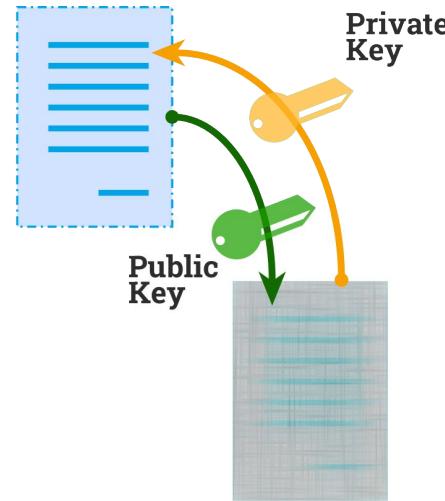




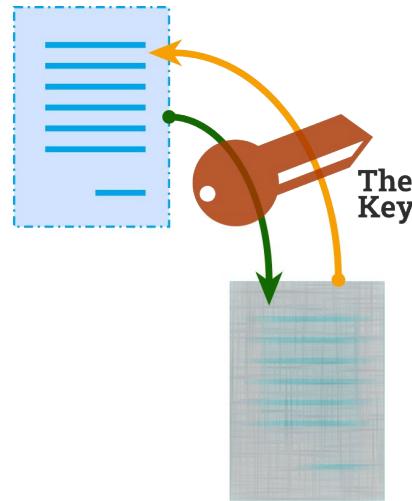
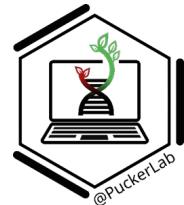
Symmetric encryption



Symmetric encryption



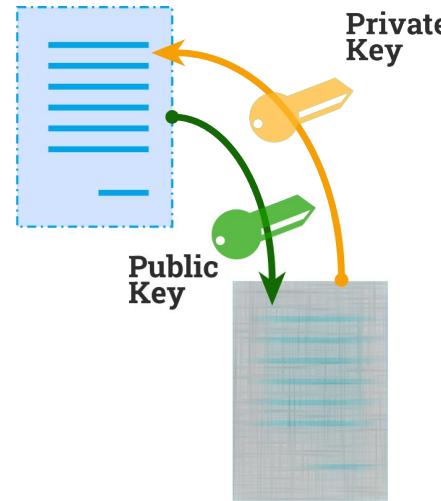
Asymmetric encryption



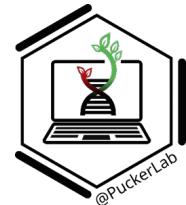
Symmetric encryption

ecdsa-sha2-nistp256

AAAAE2VjZHNhLXNoYTItbmlzdHAYNTYAAAAlbmlzdHAYNTYAAABBBJHT0aW+Q2bjkfVM7kN2ucaf8eltoalV6V0CDBgpWF+dKV1HKvGWWfjtYMOdLDUKIAzJc5Q6bKs6bkYiqqtUQ=



Asymmetric encryption



- high performance computing => server farm
- de.NBI offers VMs (VM = virtual machine)
  - running Linux
  - access through SSH
- de.NBI connected LifeScience AAI Account needed
- Tasks:
  1. Register: <https://cloud.denbi.de/wiki/registration/>
  2. Apply for the course de.NBI project:  
<https://signup.aai.lifescience-ri.eu/fed/registrar/?vo=denbi&group=pbbcourse20251>
  3. create SSH key pair in de.NBI cloud portal and confirm unix user name



- SSH key might need the correct permissions
- needed information:
  - username
  - ip address/domain
  - port number (usually optional but mandatory for de.NBI VMs)
  - ssh private key
  - exemplary commands:

```
ssh -i /path/to/key username@ip-address -p port-number
```

```
ssh -i C:\Users\User\de.NBI_course\key.txt ubuntu@134.176.27.78 -p 23457
```

```
sam@dhcp068: ~$ ssh -i /home/TUBS_NC/13.DenBiCloud_access_scripts/meckoni_ecdsa.txt ubuntu@134.176.27.78 -p 30235
[...]
de.NBI cloud
http://cloud.denbi.de
cloud@denbi.de

Welcome to Ubuntu 24.04.3 LTS (GNU/Linux 6.8.0-87-generic x86_64)

System information as of Sun Dec  7 14:28:58 UTC 2025

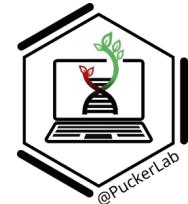
System load:  0.14      Processes:          422
Usage of /:   13.5% of 18.33GB   Users logged in:     0
Memory usage: 0%
Swap usage:   0%
IPv4 address for ens3: 192.168.0.235

Last login: Sun Dec  7 14:18:21 2025 from 192.168.0.143
ubuntu@snmtestdeletesoon-9b08e: ~$
```

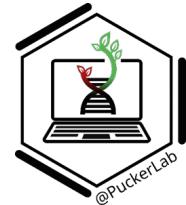


[https://github.com/bpucker/teaching  
/tree/master/deNBI2025](https://github.com/bpucker/teaching/tree/master/deNBI2025)





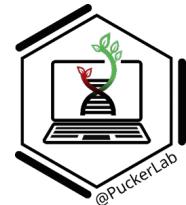
- Good structural annotation needed before functional annotation



- Good structural annotation needed before functional annotation
- BUSCO = Benchmarking Universal Single-Copy Orthologs



- Good structural annotation needed before functional annotation
- BUSCO = Benchmarking Universal Single-Copy Orthologs
- Universal genes -> expected to be in the genome
- Single-Copy -> expected to be present only once per haplotype



- Good structural annotation needed before functional annotation
- BUSCO = Benchmarking Universal Single-Copy Orthologs
- Universal genes -> expected to be in the genome
- Single-Copy -> expected to be present only once per haplophase
- if all BUSCO genes are present in the assembly => good assembly
- some might be natively missing due to gene loss in evolution

- exemplary commands:

list datasets:

```
sudo docker run --rm -v /vol/data:/vol/data ezlabgva/busco:v6.0.0_cv1 busco --list-datasets
```

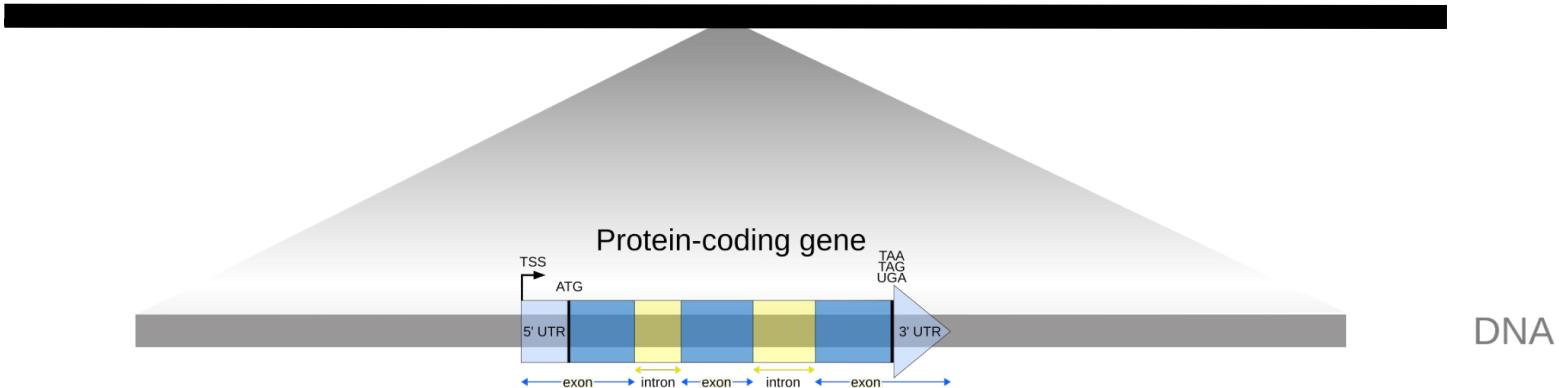
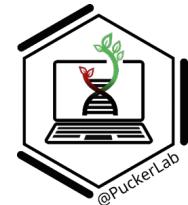
docker container name (will be automatically downloaded)

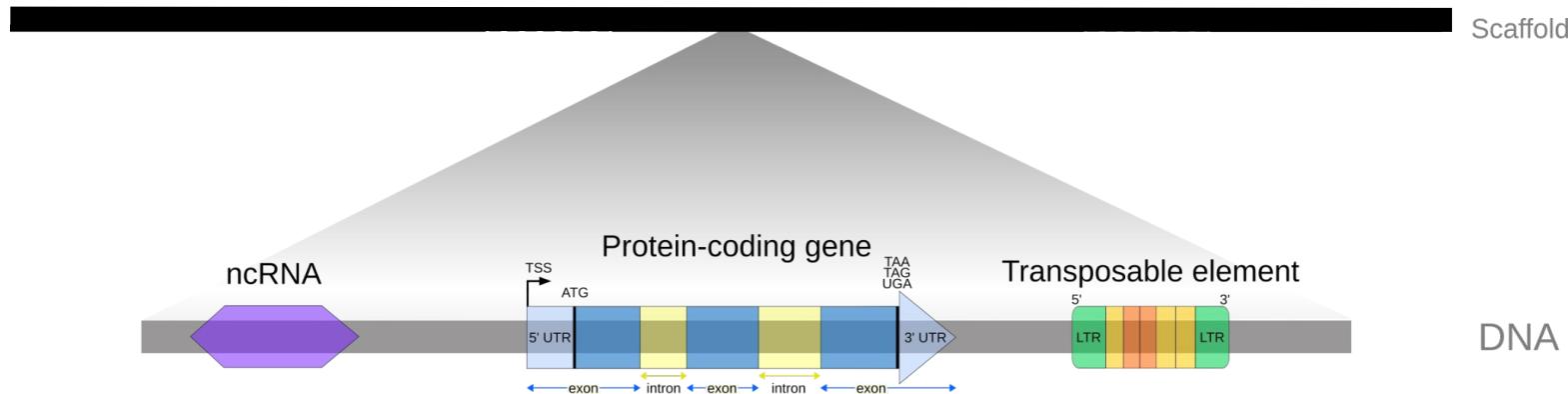
command that will be run in the  
docker container starts here

start the docker container  
input file  
mode and lineage-dataset  
output file

important for the docker container to access files

```
sudo docker run --rm -v /vol/data:/vol/data ezlabgva/busco:v6.0.0_cv1 \  
busco -i /vol/data/user/protein.fasta \  
-m proteins --cpu 13 -l brassicales_odb12 \  
--out_path /vol/data/user/BUSCO_pep/ -o busco_run_id --opt-out-run-stats
```



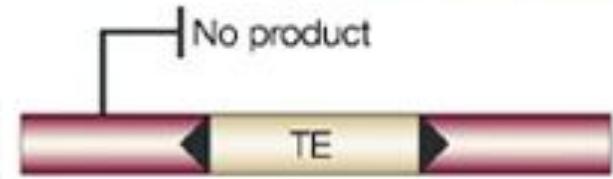
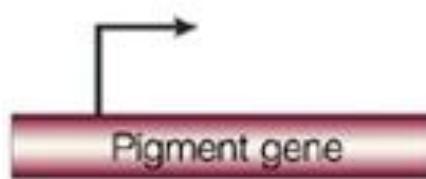


Not just protein-coding genes, but also non-coding DNA like  
Transposable element (TEs) and non-coding RNAs (ncRNAs)

In plants, only ~1-10% protein-coding genes



- Discovered by Barbara McClintock in maize
- TEs responsible for large plant genome sizes
- Based on where they land, they dictate the regulation and impact



# Classes of TEs

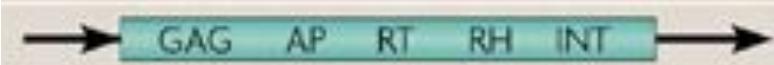


Classification	Structure	TSD	Code	Occurrence
Order	Superfamily			
<b>Class I (retrotransposons)</b>				
LTR	Copia	→ GAG AP INT RT RH →	4–6	RLC P,M,F,O
	Gypsy	→ GAG AP RT RH INT →	4–6	RLG P,M,F,O
	Bel-Pao	→ GAG AP RT RH INT →	4–6	RLB M
	Retrovirus	→ GAG AP RT RH INT ENV →	4–6	RLR M
	ERV	→ GAG AP RT RH INT ENV →	4–6	RLE M
DIRS	DIRS	→ GAG AP RT RH YR ←	0	RYD P,M,F,O
	Ngoro	→ GAG AP RT RH YR → → >	0	RYN M,F
	VIPER	→ GAG AP RT RH YR → → >	0	RYV O
PLE	Penelope	← RT EN →	Variable	RPP P,M,F,O
LINE	R2	RT EN	Variable	RIR M
	RTE	APE RT	Variable	RIT M
	Jockey	ORFI APE RT	Variable	RIJ M
	L1	ORFI APE RT	Variable	RIL P,M,F,O
	I	ORFI APE RT RH	Variable	RII P,M,F
SINE	tRNA	—	Variable	RST P,M,F
	7SL	—	Variable	RSL P,M,F
	5S	—	Variable	RSS M,O
<b>Class II (DNA transposons) - Subclass 1</b>				
TIR	Tc1-Mariner	→ Tase* →	TA	DTT P,M,F,O
	hAT	→ Tase* →	8	DTA P,M,F,O
	Mutator	→ Tase* →	9–11	DTM P,M,F,O
	Merlin	→ Tase* →	8–9	DTE M,O
	Transib	→ Tase* →	5	DTR M,F
	P	→ Tase →	8	DTP P,M
	PiggyBac	→ Tase →	TTAA	DTB M,O
	PIF-Harbinger	→ Tase* ORF2 →	3	DTH P,M,F,O
	CACTA	→ Tase ORF2 →	2–3	DTC P,M,F
Crypton	Crypton	→ YR →	0	DYC F
<b>Class II (DNA transposons) - Subclass 2</b>				
Helitron	Helitron	RPA Y2 HEL	0	DHH P,M,F

LTR copia

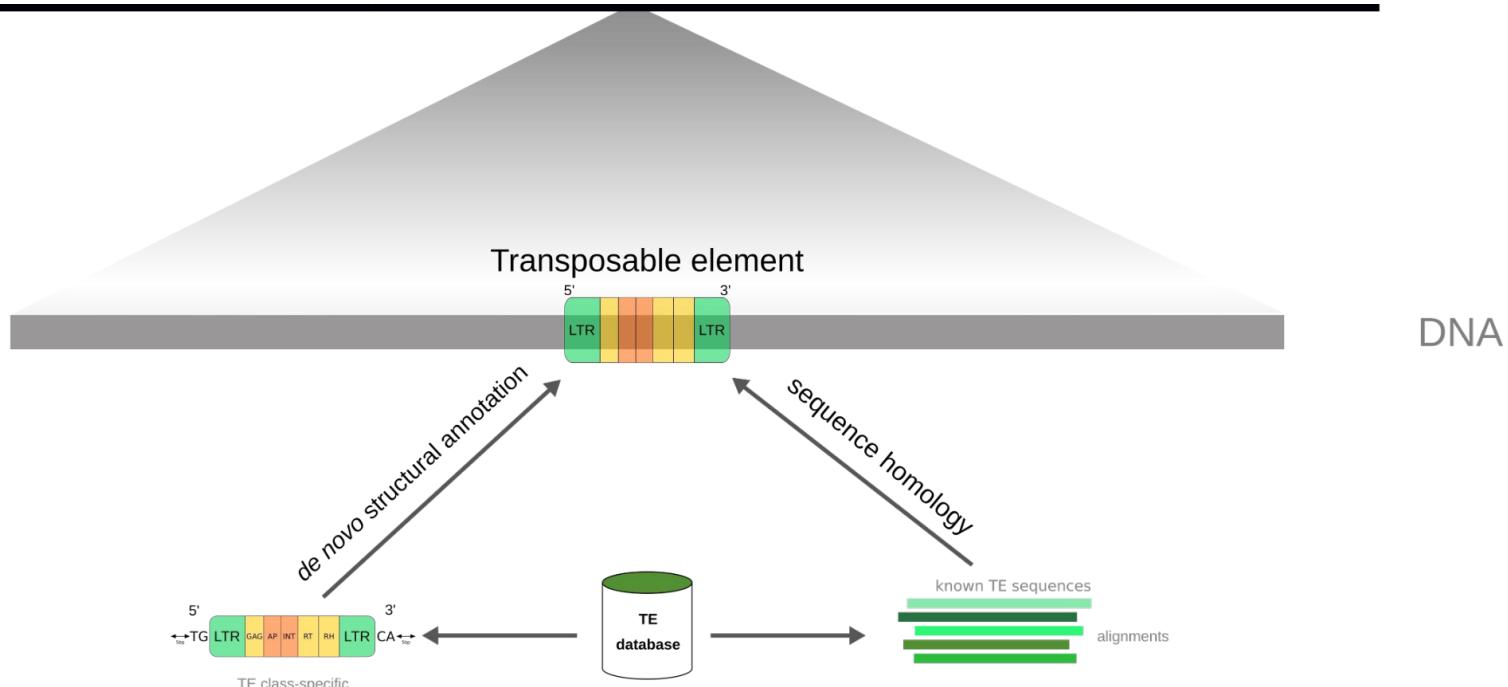
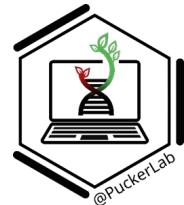


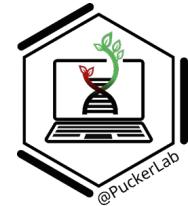
LTR gypsy



LINE L1







## Specific tools

- LTR\_FINDER/ LTR\_HARVEST → LTRs
- TIR-LEARNER → TIRs
- Helitron scanner → Helitrons

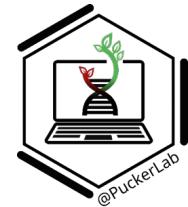
## Specific tools

- LTR\_FINDER/ LTR\_HARVEST → LTRs
- TIR-LEARNER → TIRs
- Helitron scanner → Helitrons

## EDTA (Extensive *de novo* TE annotator)

- pipeline combining above tools
- appropriate filtering
- Final TE library





### Installing Miniconda

```
cd /vol/data/tools ← #Important not in home
```

```
wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

```
bash Miniconda3-latest-Linux-x86_64.sh
```

```
# the above command will interactively ask 4 questions, first press ENTER to continue
```

```
# Do you accept license terms? [yes|no] → yes
```

```
#Miniconda3 will be installed in /home/ubuntu/miniconda3 or specify a different location → /vol/data/tools/Miniconda3
```

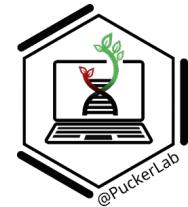
```
#Proceed with initialization [yes|no] → yes
```

```
source ~/.bashrc
```

```
conda config --add channels conda-forge
```

```
conda update -n base --all
```

```
conda install -n base mamba
```



## Installing EDTA through conda/mamba

```
mamba create -n EDTA2.2 -c conda-forge -c bioconda -c r annosine2 biopython cd-hit coreutils genericrepeatfinder  
genometools-genometools glob2 tir-learner ltr_finder_parallel ltr_retriever mdust multiprocess muscle openjdk perl perl-text-soundex  
r-base r-dplyr regex repeatmodeler r-ggplot2 r-here r-tidyr tesorter samtools bedtools LTR_HARVEST_parallel HelitronScanner
```

## Getting the latest version of EDTA from GitHub

```
cd /vol/data/tools; git clone https://github.com/oushujun/EDTA.git; cd EDTA; git checkout EDTA2; git branch
```

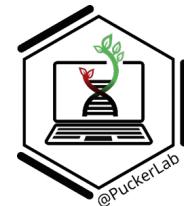
## Running EDTA

```
cd /vol/data/user
```

```
mkdir EDTA_results; cd EDTA_results
```

```
conda activate EDTA2.2
```

```
perl /vol/data/tools/EDTA/EDTA.pl --genome /vol/data/DatasetB.genome.fasta --overwrite 1 --anno 1 --sensitive 1 --evaluate 1 --threads 13  
>>edta.log 2>&1 &
```



**DatasetB.genome.fasta**

DatasetB.genome.fasta.mod

DatasetB.genome.fasta.mod.EDTA.TEanno.density\_plots.pdf

DatasetB.genome.fasta.mod.EDTA.TEanno.gff3

DatasetB.genome.fasta.mod.EDTA.TEanno.sum

DatasetB.genome.fasta.mod.EDTA.TElib.fa



**DatasetB.genome.fasta.mod.EDTA.anno**

**DatasetB.genome.fasta.mod.EDTA.combine**

**DatasetB.genome.fasta.mod.EDTA.final**

DatasetB.genome.fasta.mod.EDTA.intact.fa

DatasetB.genome.fasta.mod.EDTA.intact.gff3

**DatasetB.genome.fasta.mod.EDTA.raw**

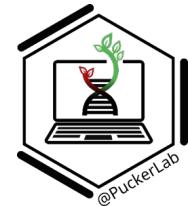
**DatasetB.genome.fasta.mod.MAKER.masked**

**edta.log**

**DatasetB.genome.fasta.mod.RM2.raw.fa**

**edta2.log**

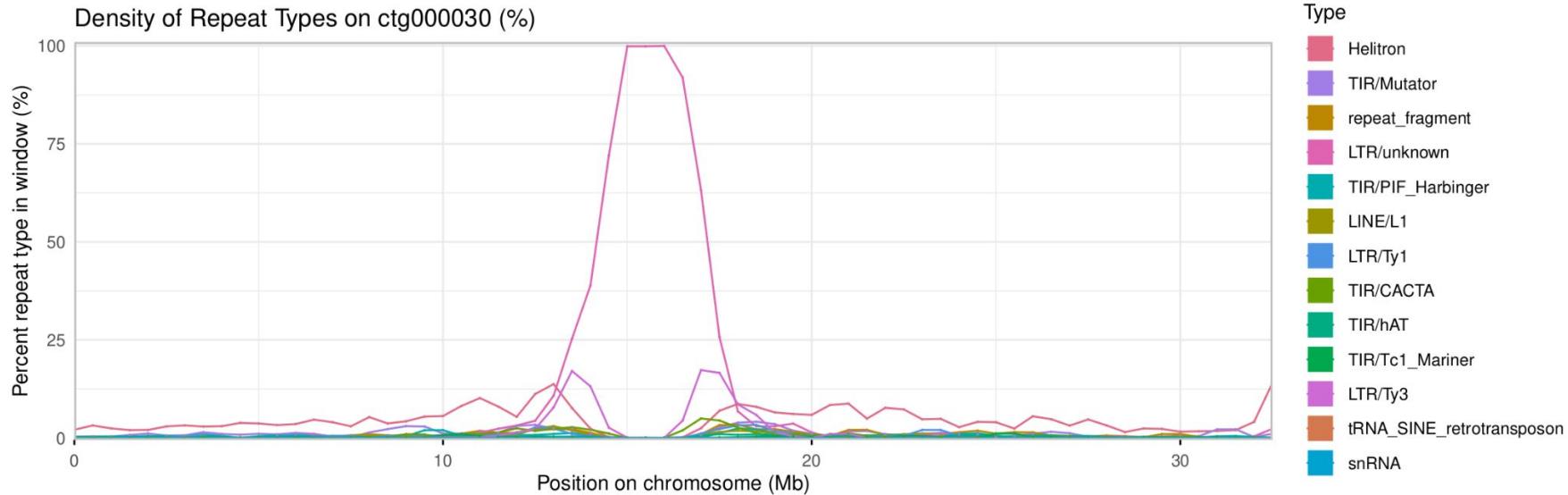
**DatasetB.genome.fasta.mod\_divergence\_plot.pdf**



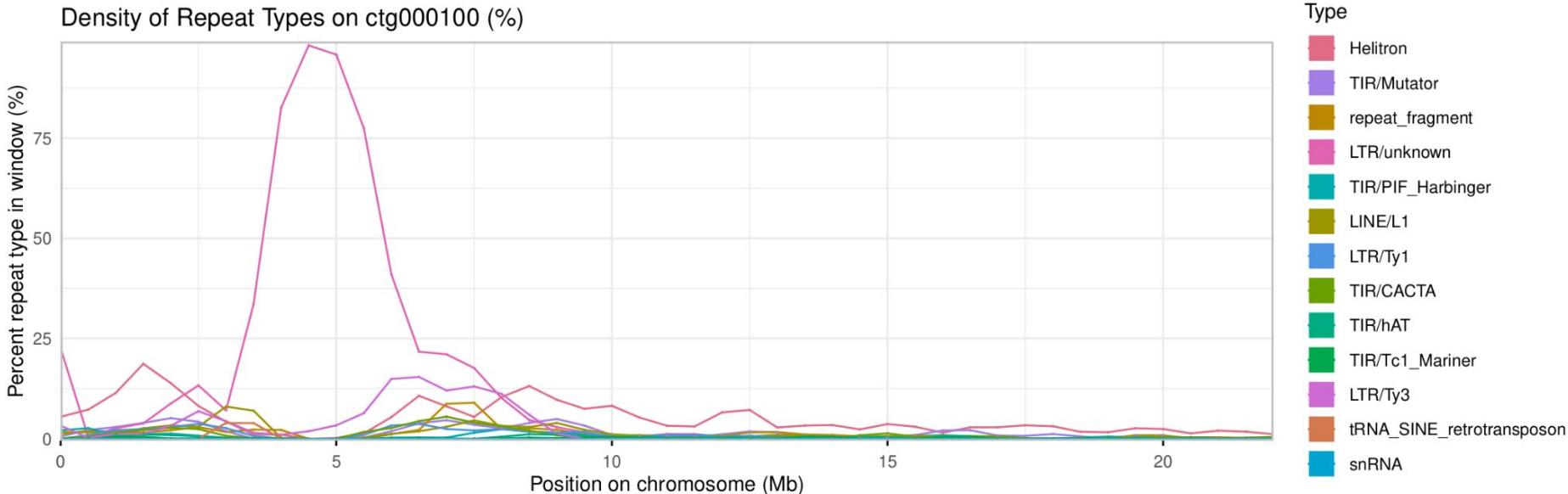
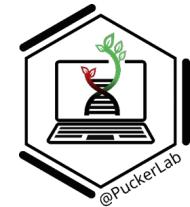
1. \*.EDTA.TEanno.sum →
2. \*.EDTA.TEanno.gff3
3. Density plots for top contigs

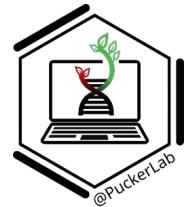
Repeat Classes			
=====	=====	=====	=====
Total Sequences: 13			
Total Length: 134799993 bp			
Class	Count	bpMasked	%masked
=====	=====	=====	=====
LINE	--	--	--
L1	1239	989157	0.73%
LTR	--	--	--
Copia	738	1002959	0.74%
Gypsy	2081	3195026	2.37%
unknown	5660	16713457	12.40%
SINE	--	--	--
tRNA	143	179011	0.13%
TIR	--	--	--
CACTA	1668	918659	0.68%
Mutator	3170	1767375	1.31%
PIF_Harbinger	1751	728452	0.54%
Tc1_Mariner	539	259836	0.19%
hAT	807	335497	0.25%
nonTIR	--	--	--
helitron	11058	5960229	4.42%
repeat_fragment	3879	1580072	1.17%
-----			
total interspersed	32733	33629730	24.95%
-----			
snRNA	8	859	0.00%
-----			
Total	32741	33630589	24.95%

## Density plot



# Density plot





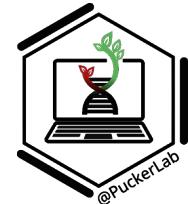
### Today you learned:

- Accessing de.NBI VMs via SSH
- Working in a CLI environment
- Running BUSCO
- Identifying transposon elements

### Tomorrow you will learn:

- tRNA annotation
- Running InterProScan
- KIPEs3, MYB\_annotator
- Tree building

## Availability of materials

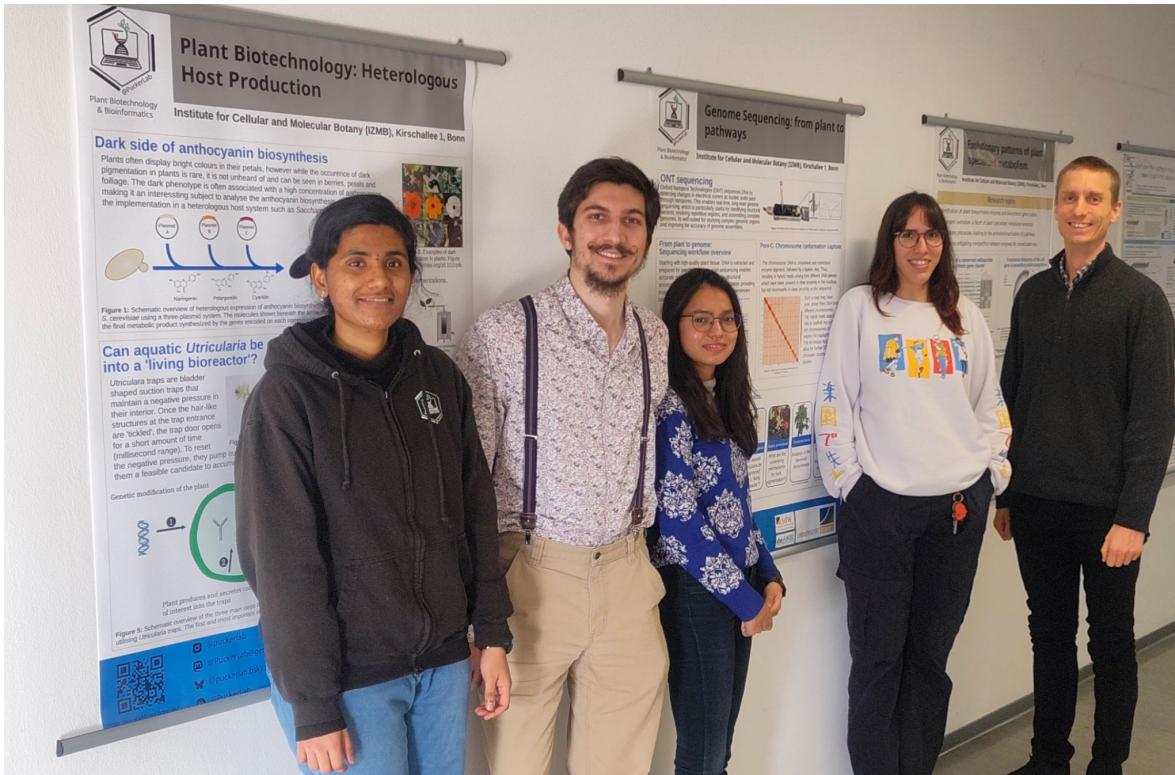


Please find all course materials here:

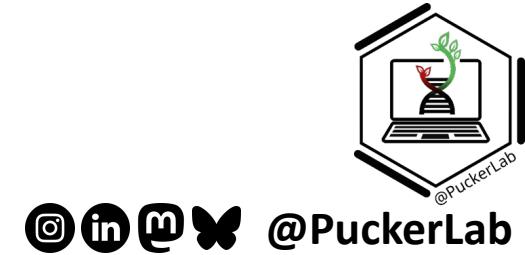
<https://github.com/bpucker/teaching/tree/master/deNBI2025>



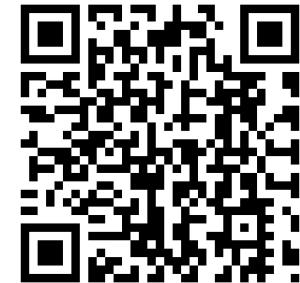
# Acknowledgements



From Gene Models to Biological Insights - Functional Annotation | Day 1 | 34



Email: [pucker@uni-bonn.de](mailto:pucker@uni-bonn.de)  
 Web: [pbb.uni-bonn.de](http://pbb.uni-bonn.de)



Supported by:

