

***Busseola fusca* Genome Consortium Overview**

**Using Bioinformatics and
Genomics Tools to
Develop Integrated, Reproducible,
Shareable Workflows**

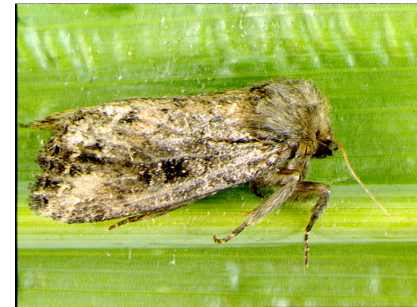
Agenda

- *Busseola fusca* Genome Consortium
(tea break)
- The importance of reproducibility and collaboration
 - Data cleaning/ Software version/Parameter choice
 - Record keeping
 - Sharing workflows
- Hands on exercise using Galaxy

Sequencing the Whole Genome and Transcriptome of *Busseola fusca* here at ILRI



Busseola fusca is a major crop pest, especially of maize



B. Fusca devastates large proportions of both small & large-scale crops in sub-Saharan Africa

TABLE 1. Mean percentages of plants infested with eggs (\pm SE), number of egg batches per plant, discovery efficiency and percentage egg parasitism on *Busseola fusca* at the whorl stage of maize in four agroecological zones (AEZ) in Uganda during 2004^z

Season	AEZ	Egg batches/plant	Discovery efficiency (%)	% Egg parasitism
1st rains 2004	Eastern	0.010 \pm 0.005 a	9.8 \pm 6.0 a	39.4 \pm 16.9 a
	Southeastern	0.016 \pm 0.005 a	9.6 \pm 4.8 a	17.7 \pm 8.4 a
	Lake Albert Crescent	0.00 b	NF	NF
	Lake Victoria Crescent	–	–	–
	<i>df</i>	2, 62	1, 47	1, 19
	<i>F</i> value	4.11	1.05	1.32
	<i>P</i> value	0.0211	0.3465	0.2642
	Eastern	0.007 \pm 0.003 a	5.6 \pm 5.6 ab	27.2 \pm 12.4 ab
	Southeastern	0.020 \pm 0.005 a	18.8 \pm 7.4 a	46.1 \pm 10.2 a
	L. Albert Crescent	0.015 \pm 0.006 a	1.8 \pm 1.8 ab	7.4 \pm 5.0 b
L. Victoria Crescent	0.004 \pm 0.002 a	0.0 \pm 0.0 a	0.0 \pm 0.0 b	
<i>df</i>	3, 73	3, 73	3, 30	
<i>F</i> value	2.48	2.99	3.90	
<i>P</i> value	0.0680	0.0365	0.0183	

means followed by a common letter do not differ significantly at $P \leq 0.05$ (SNK).

ed.

NF, host not found.

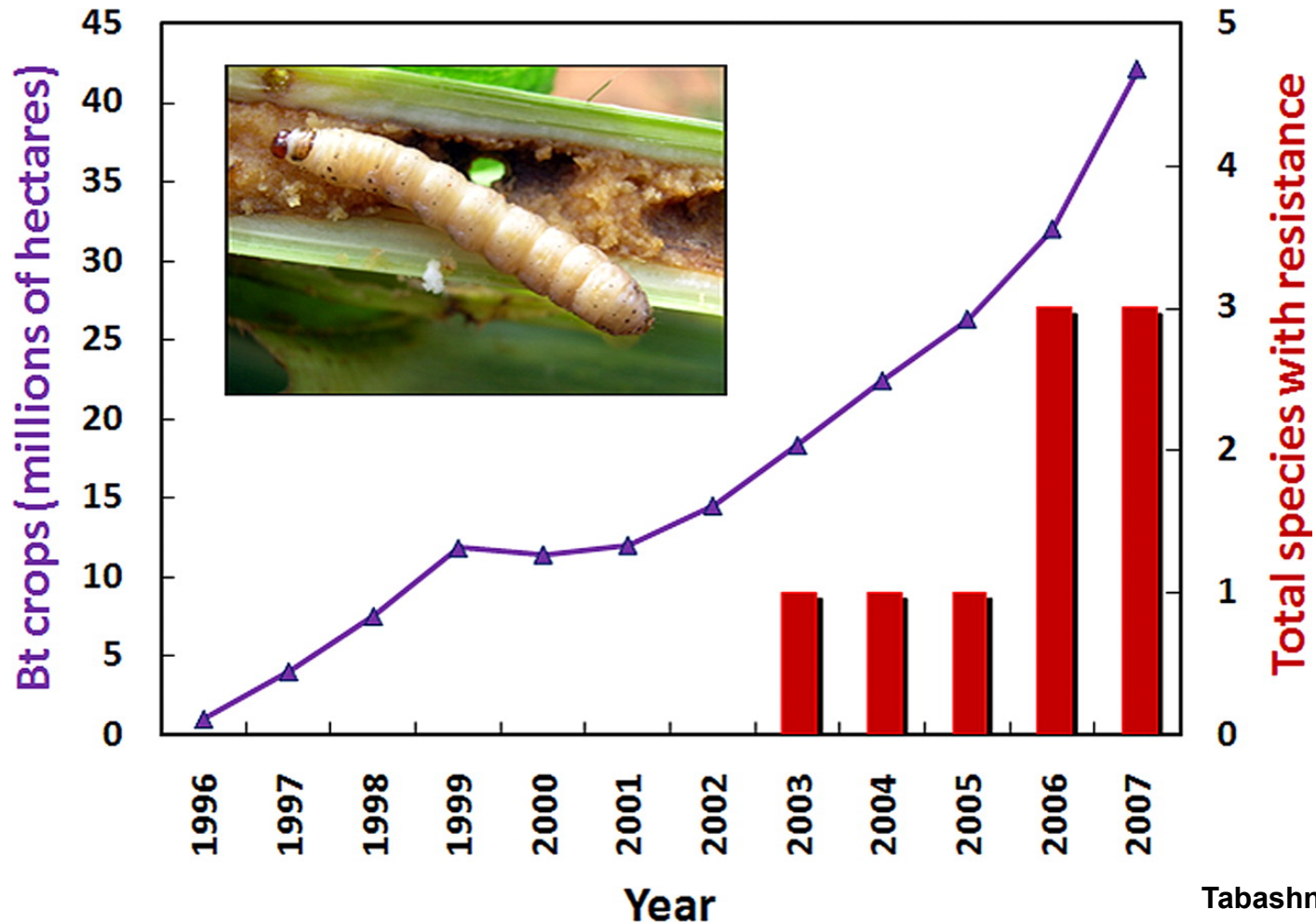


Current Methods of Pest Management

Cultural control – destroying crop residues, intercropping, rotation, etc.

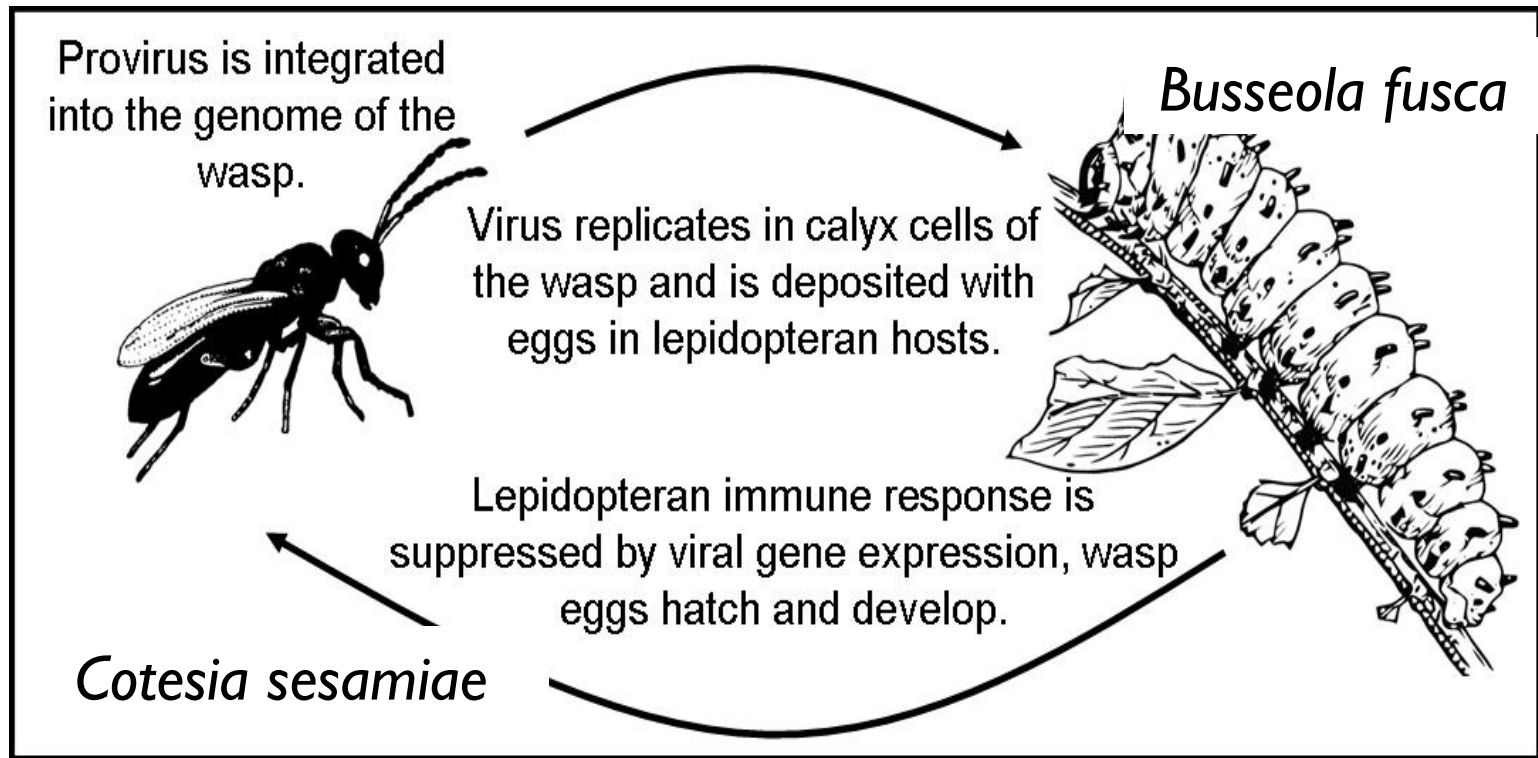
Synthetic pheromone traps – disrupt mating

Breed resistant crops – but larvae evolve resistance to the resistance!



Biological Control for Pest Management

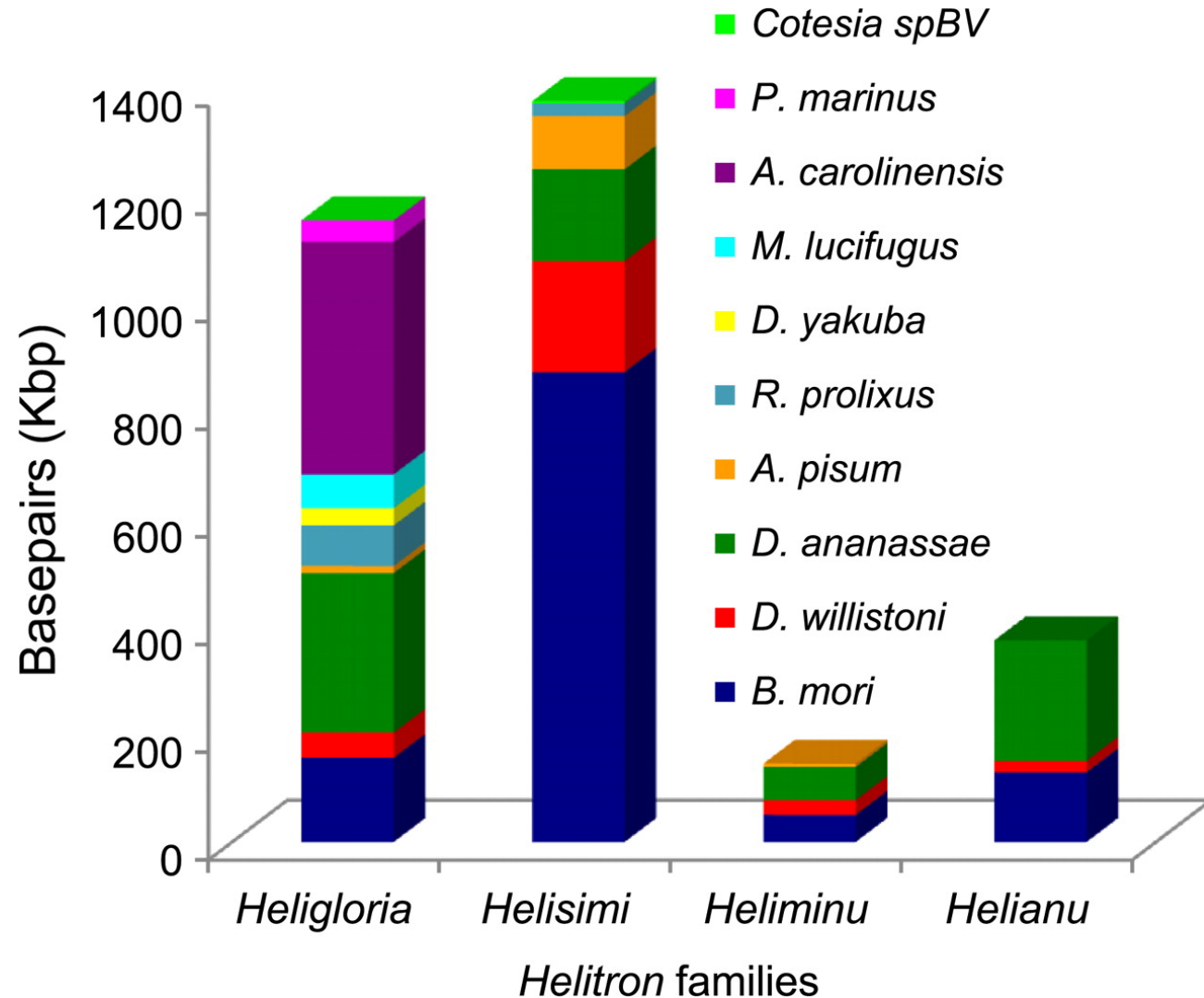
Parasitoid wasps use *B. fusca* larvae to lay eggs, and simultaneously deposit viral particles (the virus is embedded in the wasp's genome).



Fewer chemicals, affordable, widespread/sustainable, efficient.

Multiple Reasons to Sequence *B. fusca*

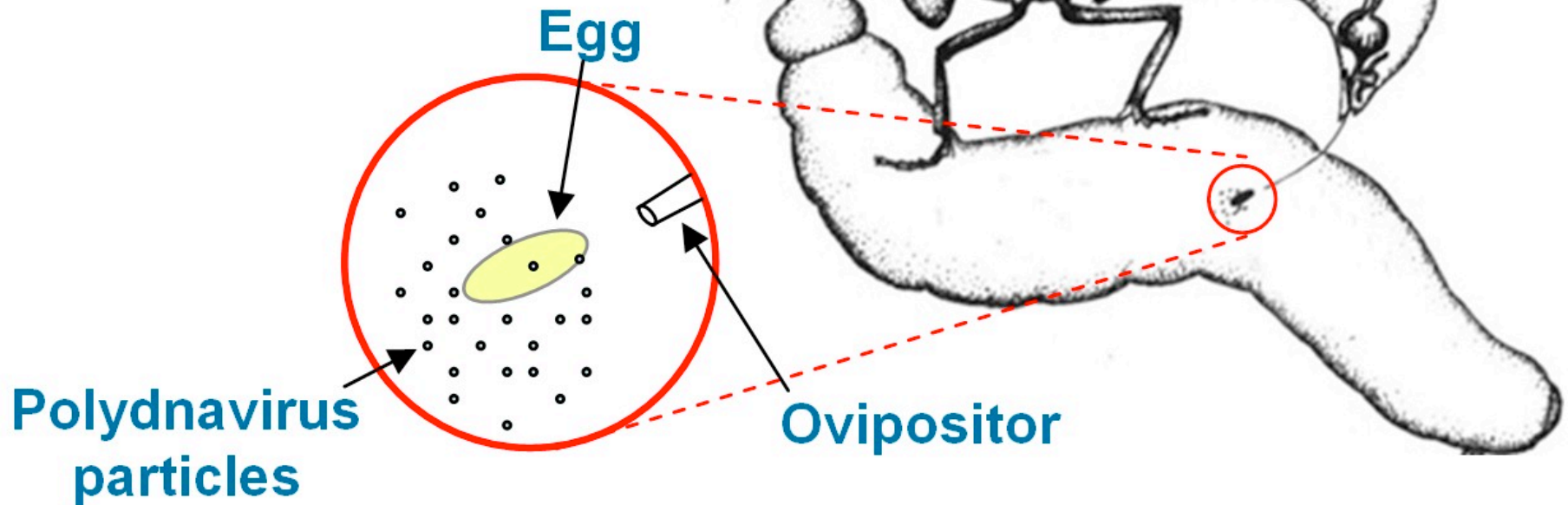
- **Major crop pest threatening food security in East Africa**
- Previous cases of HTT involve the virus in the parasitoid wasp of *B. fusca*



Closely-Associated Species Can Exchange DNA Horizontally



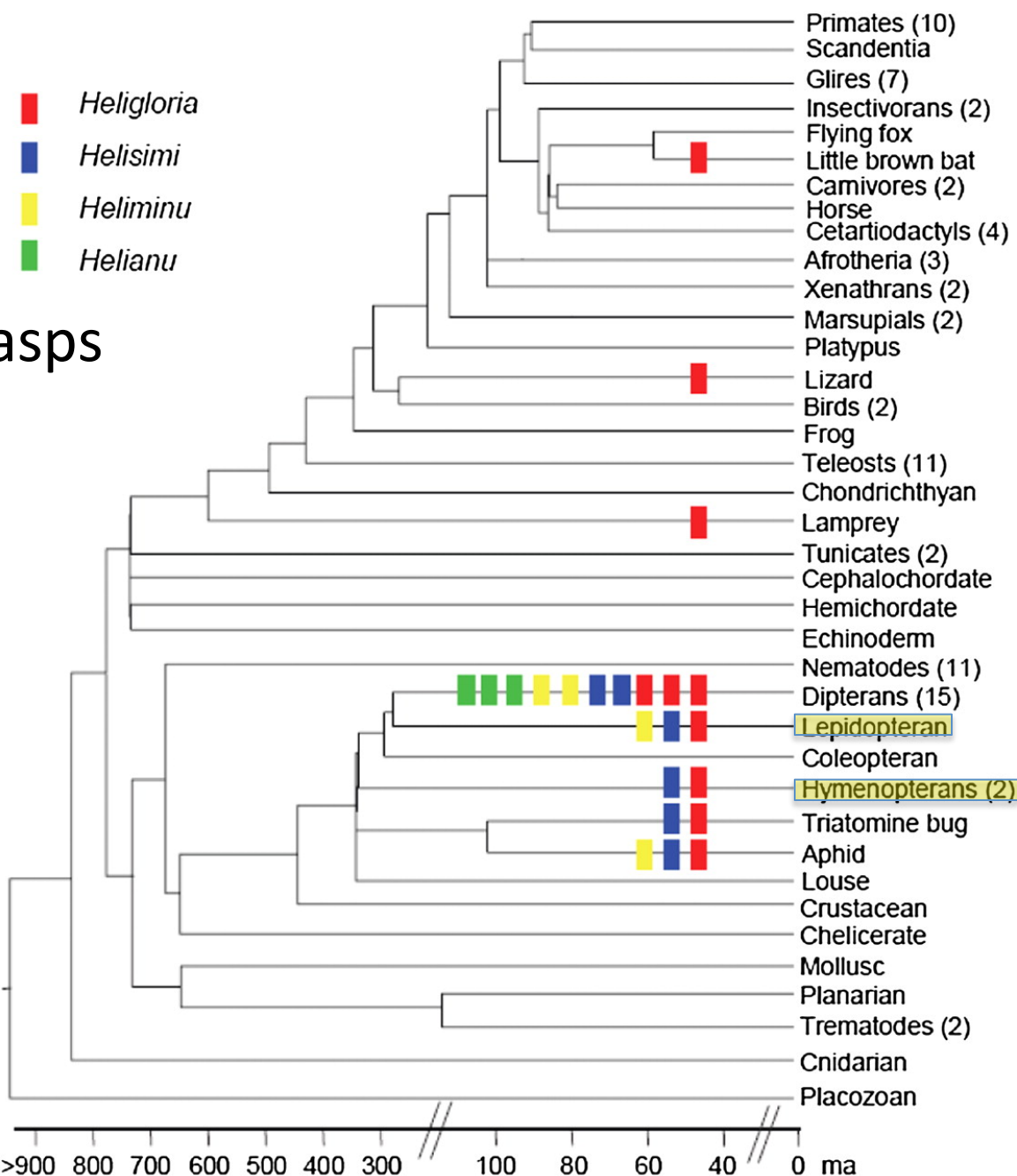
The virus suppresses the immune response, thus the “infection” usually, but not always, kills the larvae.



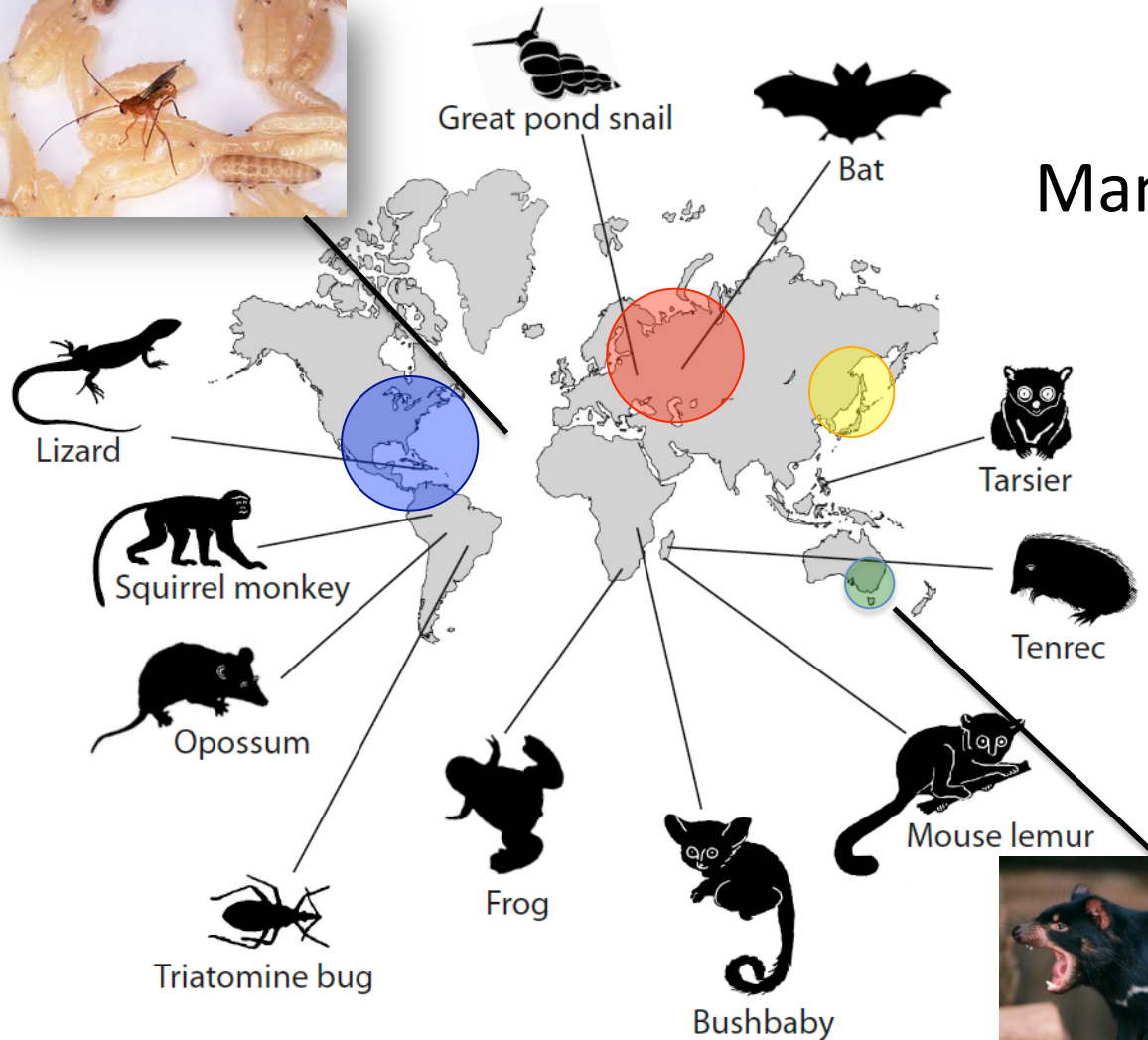
Nearly identical *transposable elements* have been found in the WGS of various organisms including:

the viruses in *Cotesia* wasps and the first sequenced lepidopteran, *Bombyx mori*

A TE is a piece of DNA that is, or once was, capable of moving or replicating and reinserting in the genome.



Unexpectedly, *many* cases of horizontal transfer among eukaryotes found in WGS data



Many different types of TEs

Many species (>13)

Spanning 4 phyla

On 4 continents

But how?

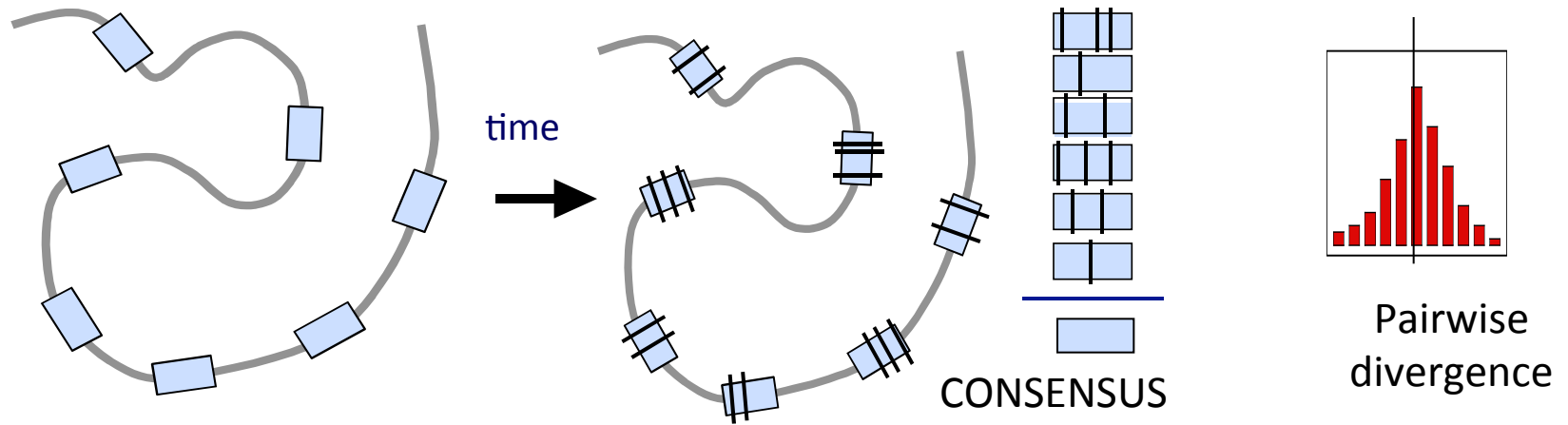


Gilbert/Schaack et al. 2010
Schaack/Gilbert et al. 2010
Gilbert et al. and Schaack 2013

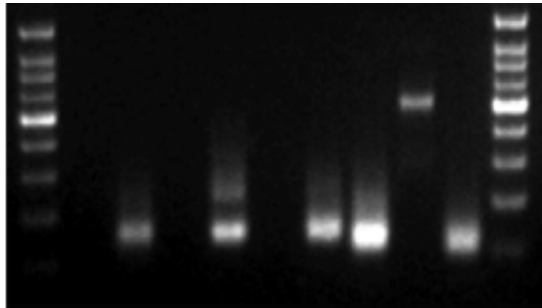
Searching WGS Data for Evidence of HTT

```
>contig00001 length=577 numreads=6
TGC GGGACTCGGTAGACGAATTGCTCAGGTGCTTCTCGCTCCTGCTCGGAAaCTttGTTTC
TCAGCCTCGTAATTTCTCTGATTCTGGTACTTGTTATCGTTCGAACTGCGCCGGGCCGTAG
TTTTTCGTCGGCTGAAGACTTGACTTGCTGGTGATGATTGTTGGGGTGGTAGACCAGCTTC
GGGCGTTCCATCCTGACGCGTTGGATAATTTTCGGGTTtGTTtCGTCGTGGTGAACAGCGGA
CGAGGAGTCGAAAGTGACCCTGCGTGGTAGACGTCAGAGTGTTTCGTCGCTGCTTTCTTGG
GATTCGCGAACCTTGATGGTAGTTTCGTCTTGCAGCTGCCGGCGTAAGTTATTGCCACTG
ACCGGCGACTTGGACTCTTCTtCTTTTTCTCAAACCTTTCAGTTATTGGACGGAGGGTG
TGATCAAGCGGCCGTTTTTGTGAAGTAGCTCTTACGCTGCGTGGCGGCTTTGGCTTCGACG
GTTGGCCGAGTGATCTGCGAGTCTTCTTCCAAGTTGTCCGGCGATTCTTGCTCACTGGAC
TCTGGCTCCGGGCTTGATCTCaagTTTGTGAACAaTC
>contig00002 length=3271 numreads=157
aaTagAatCtAATTTTTATATAAaTCACTTTAttCTTAttATTAaCTTGATGAAAaCTTT
TTtCTTCGTGTTAAAACATATATTTAGATTCTGACCTGTTATTGTTTTCAAACCTTAGTTT
GGCTTCAGCTTACTTATAGATGCTTTAAGTTACCTACTTCATTTTAGTGTGACTCCACCG
TAGATATCGAGCCAAAAATTCGTATCTGTGTAACGCTTAAATAGAAACATCAAGTCGAAT
TTGTCGATATCTATTCATCTATCAGTAAGTCTTACAACGTTGGATTCACTTCTGTTCTTATT
TTCAACAAAAAAGTGTATTTTTTTTTtCAaTGGTCCGAGCGAATCCGCAAGAAGTAATGCAG
TCTGCATAACGTTGAATTTATGAGCAACTAACCTGCTCTGTTTTTGCAAGGACTAGAATCA
GTTGCATTGAGCAACATTGCTTTCAAGTCAGTAGAGTCTTCTTGCAATATAAGCGCTAAA
AAAGTTTTCTTGAACAACATTAGAGTTCCGAGTTCCCAACAAGCTTTCTAAATAACTGGA
ACGGAAGTTGTTCTAACTAACAAACGAATTCGTCGGACTGGAACAGGAACATTACATAGTG
GGCCATAAAGTG CATCTAGAAAATAATACATATCACGGCCAACCTCCAACCTTCATCAAACA
GCAGGAAATGAAATTGACATAGTAAGAAATATCTATTATGGTCAAGTCCAAATCCATGAG
T TACTCGTAACTCAAGTAAAACgaTCGTTCAATACTTATAATAAACTTCCAATAACTTA
CCTATTGGAAACTGTGCAAGATTTATTCAAAGTGTGTATAATAAATCGAAGGCTGGAAGT
GAGGTAGCTTTAAGGCTCGTTAGAAAAATTACTCGCAAGTGCAGTACTTTTTCTTCGGAA
ATGTAATAACAACACTTAGTCGTAATTTTTTTTTTTTTtGAaTAACTTTTATTGTAaCTA
```

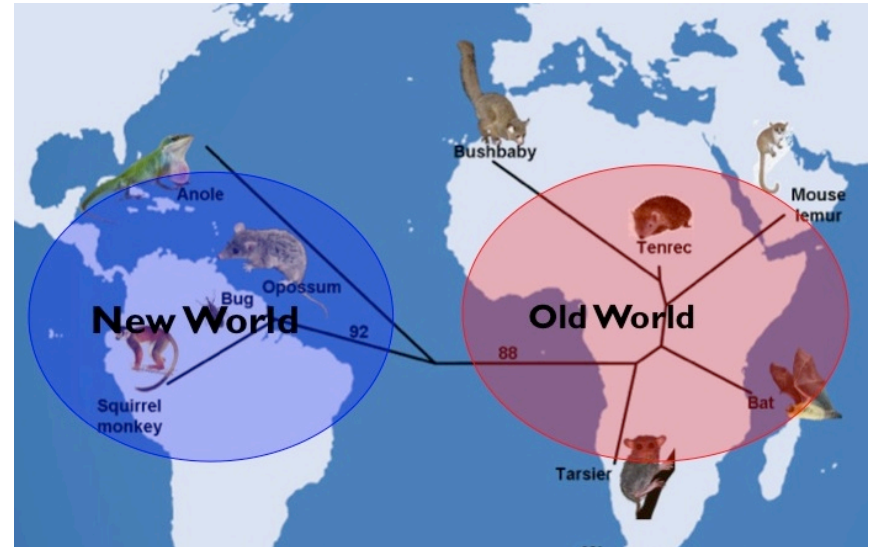
Searching WGS Data for Evidence of HTT



PCR/
Cloning/
Sequencing
to validate

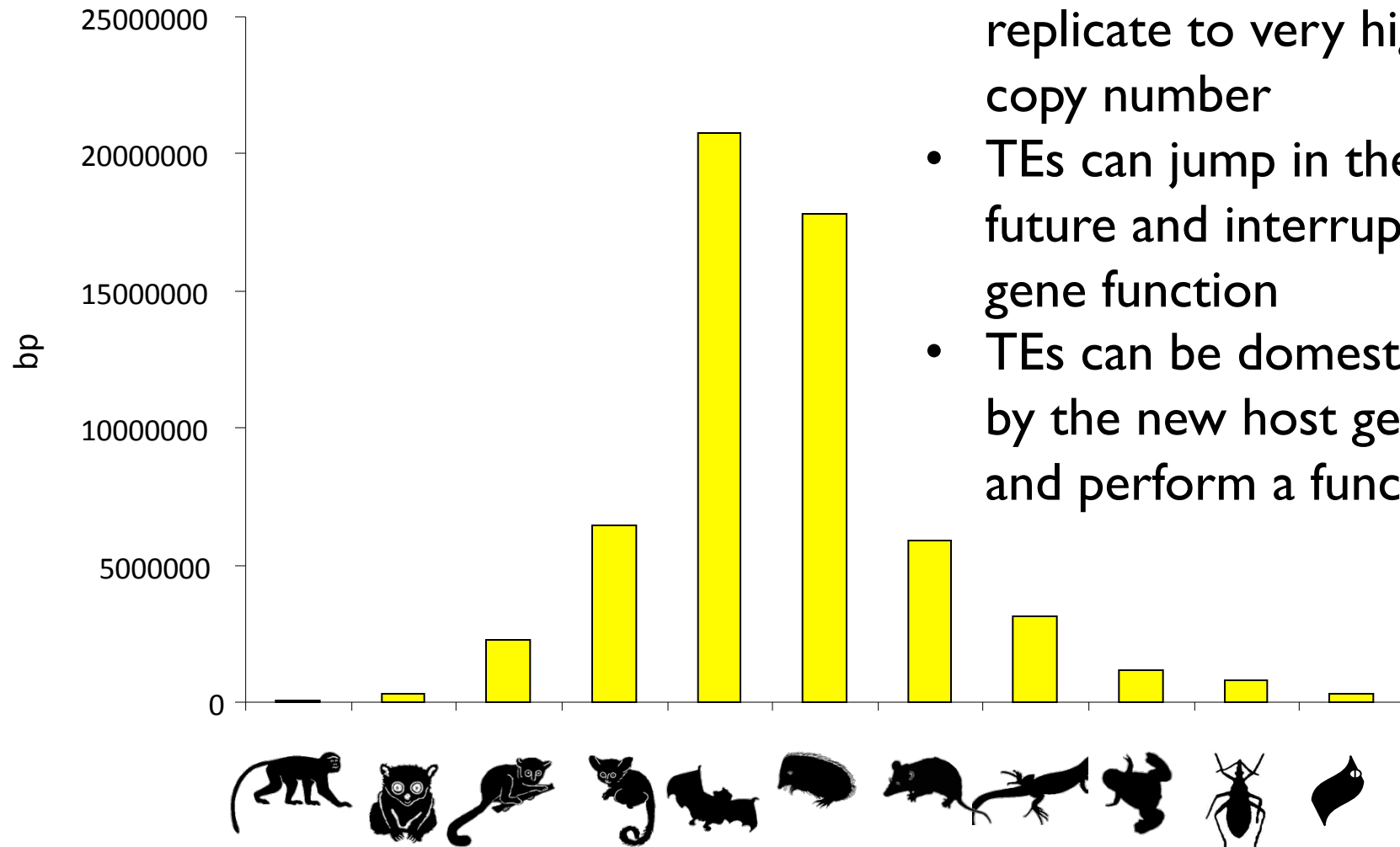


$$\text{Average divergence} \times \text{Mutation Rate} = \text{Calculate the date of invasion}$$



Biogeographical
& natural history data

But who cares if TEs are exchanged horizontally?



- After HT, TEs can replicate to very high copy number
- TEs can jump in the future and interrupt host gene function
- TEs can be domesticated by the new host genome and perform a function

In bacteria, we knew HTT was frequent and important.

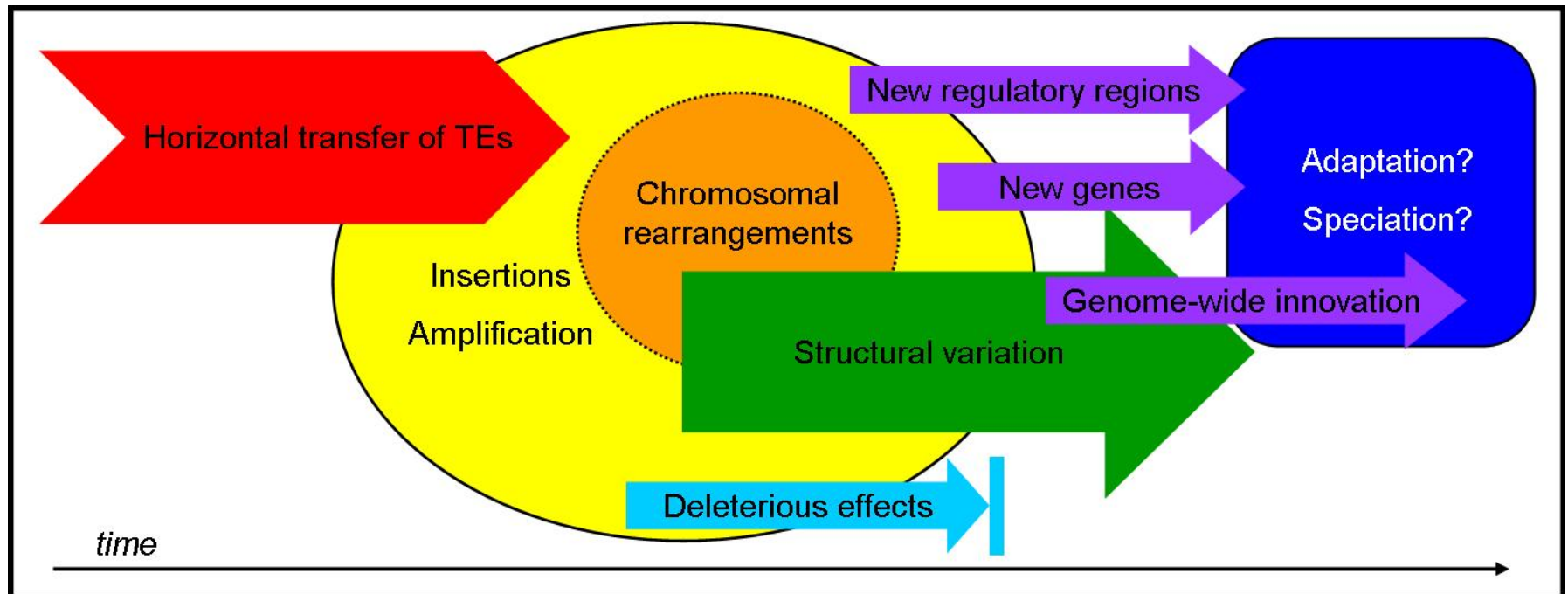
**In eukaryotes,
we now know that HTT is not uncommon.**

What are the mechanisms?

How are TEs incorporated into the germline?

What is the impact?

How often does this happen?



WGS and Transcriptome Sequencing

Genome sequencing:

DNA extracted from a single larva (1 run, ~15x coverage of 900 MB)

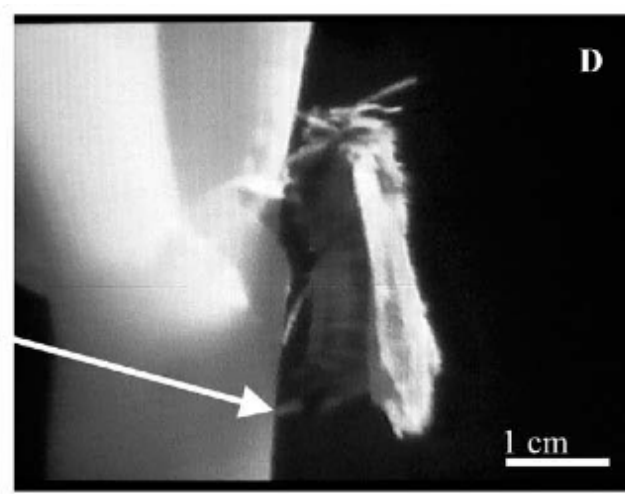
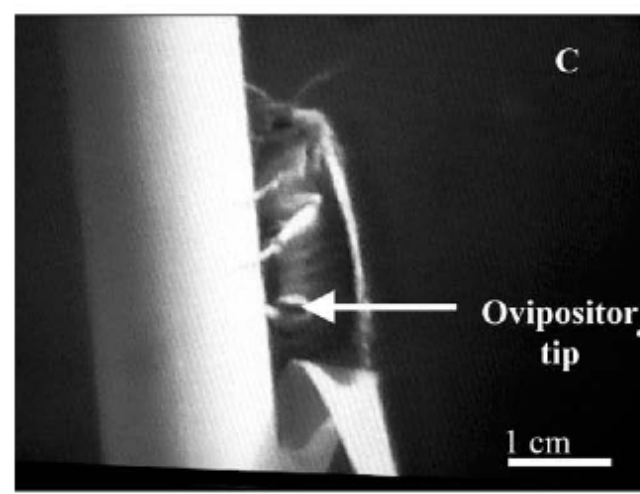
WGS and Transcriptome Sequencing

Genome sequencing:

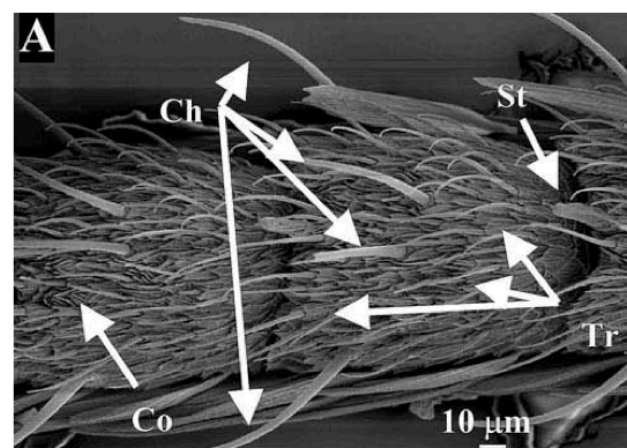
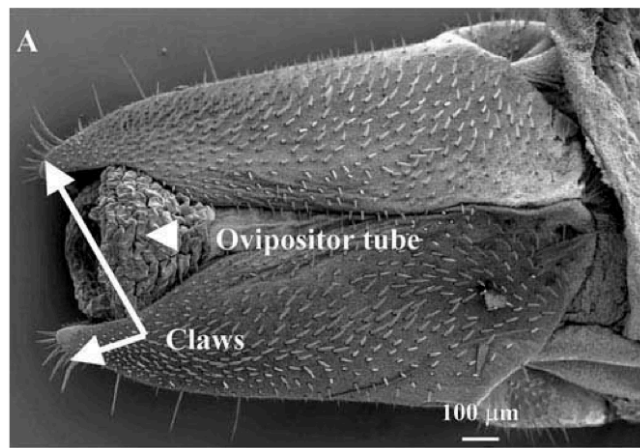
DNA extracted from a single larva (2 runs, ~30x coverage)

Transcriptome sequencing:

S2_S1_L001_R1_001.fastq	B. fusca unexposed; larvae (1)
S3_S2_L001_R1_001.fastq	B. fusca exposed to CSM; larvae (1)
S6_S3_L001_R1_001.fastq	B. fusca female thorax (1)
S7_S4_L001_R1_001.fastq	B. fusca male thorax (1)
S10_S5_L001_R1_001.fastq	B. fusca female ovipositor (1)
1_S1_L001_R1_001.fastq	B. fusca neonates (350)
2_S2_L001_R1_001.fastq	B. fusca eggs (700)
3_S3_L001_R1_001.fastq	B. fusca male antennae (60)
4_S4_L001_R1_001.fastq	B. fusca female antennae (60)
5_S5_L001_R1_001.fastq	B. fusca exposed to CSK; larvae (1)



Genes upregulated in different developmental stages, organs, males versus females, and exposed and unexposed to wasps



Differential expression of the CrV1 haemocyte inactivation-associated polydnavirus gene in the African maize stem borer *Busseola fusca* (Fuller) parasitized by two biotypes of the endoparasitoid *Cotesia sesamiae* (Cameron)

C.W. Gitau^{a,c,*}, D. Gundersen-Rindal^b, M. Pedroni^b, P.J. B. ^c, and J. Dupas^d

^aCentre of Insect Physiology and Ecology, P.O. Box 30772, Nairobi, Kenya
^bA-ARS Insect Biocontrol Laboratory, Beltsville, MD, USA
^cKenyatta University, P.O. Box 43844-00100, Nairobi, Kenya
^dCentre de Populations Génétique et Evolution, 1 av. de la Forêt, F-91198 Gif-sur-Yvette, France

Received 15 September 2006; received in revised form 11 October 2006; accepted 11 April 2007

Cotesia sesamiae Mombasa (CsM) causes an immune response and can't infect *B. fusca* successfully

Abstract
The two biotypes of the endoparasitoid of *Busseola fusca* and *Sesamia calamistis* in sub-Saharan Africa exist as two biotypes. In Kenya, the western biotype completes development in *B. fusca* larvae. However, eggs of the coastal *C. sesamiae* are encapsulated in this host and ultimately, no parasitoids emerge from parasitized *B. fusca* larvae. Both biotypes develop successfully in *S. calamistis* larvae. Encapsulation activity by *B. fusca* against eggs of the avirulent *C. sesamiae* was detectable six hours post-parasitization. The differences in encapsulation activity between virulent and avirulent strains were associated with differences in nucleotide sequences and expression of a CrV1 polydnavirus (PNDV) gene. PNDV is associated with haemocyte inactivation in the *Cotesia* spp. In *B. fusca*, CrV1 expression was faint in fat body and haemolymph samples from *B. fusca* parasitized by *C. sesamiae* Mombasa (CsM), which exhibited encapsulation activity. In *S. calamistis*, CrV1 expression was high in fat body and haemolymph samples from *B. fusca* parasitized by *C. sesamiae* Kitale (CsK), which does not exhibit encapsulation activity. The CrV1 gene sequences between virulent and avirulent wasp suggest that the differences in encapsulation activity are due to qualitative differences in CrV1-haemocyte interactions.

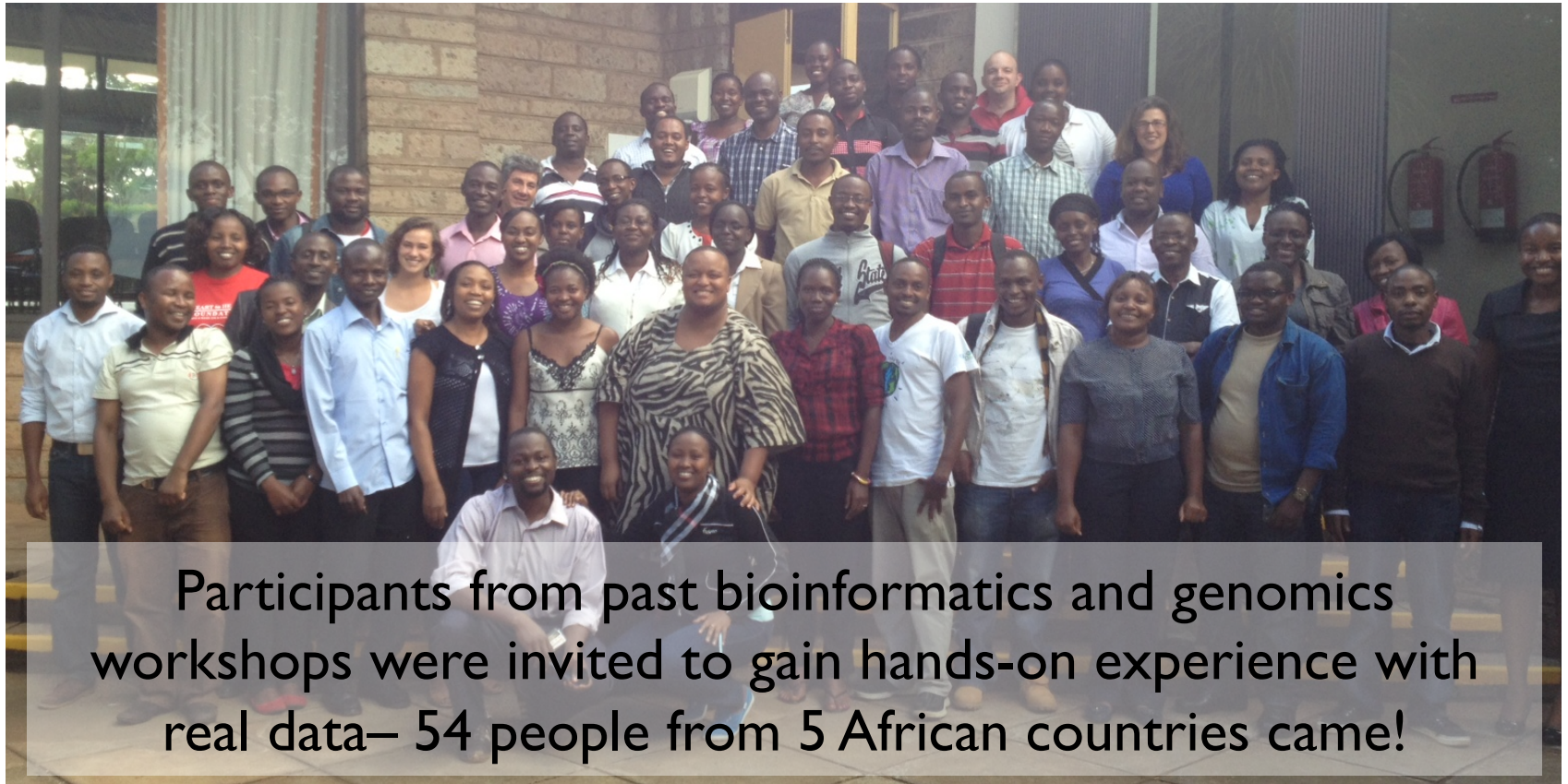
Cotesia sesamiae, a widespread endoparasitoid of *Busseola fusca* and *Sesamia calamistis* in sub-Saharan Africa exists as two biotypes. In Kenya, the western biotype completes development in *B. fusca* larvae. However, eggs of the coastal *C. sesamiae* are encapsulated in this host and ultimately, no parasitoids emerge from parasitized *B. fusca* larvae. Both biotypes develop successfully in *S. calamistis* larvae. Encapsulation activity by *B. fusca* against eggs of the avirulent *C. sesamiae* was detectable six hours post-parasitization. The differences in encapsulation activity between virulent and avirulent strains were associated with differences in nucleotide sequences and expression of a CrV1 polydnavirus (PNDV) gene. PNDV is associated with haemocyte inactivation in the *Cotesia* spp. In *B. fusca*, CrV1 expression was faint in fat body and haemolymph samples from *B. fusca* parasitized by *C. sesamiae* Mombasa (CsM), which exhibited encapsulation activity. In *S. calamistis*, CrV1 expression was high in fat body and haemolymph samples from *B. fusca* parasitized by *C. sesamiae* Kitale (CsK), which does not exhibit encapsulation activity. The CrV1 gene sequences between virulent and avirulent wasp suggest that the differences in encapsulation activity are due to qualitative differences in CrV1-haemocyte interactions.

Cotesia sesamiae Kitale (CsK) CAN infect *B. fusca*

Genomics and Transcriptomics of *B. fusca*

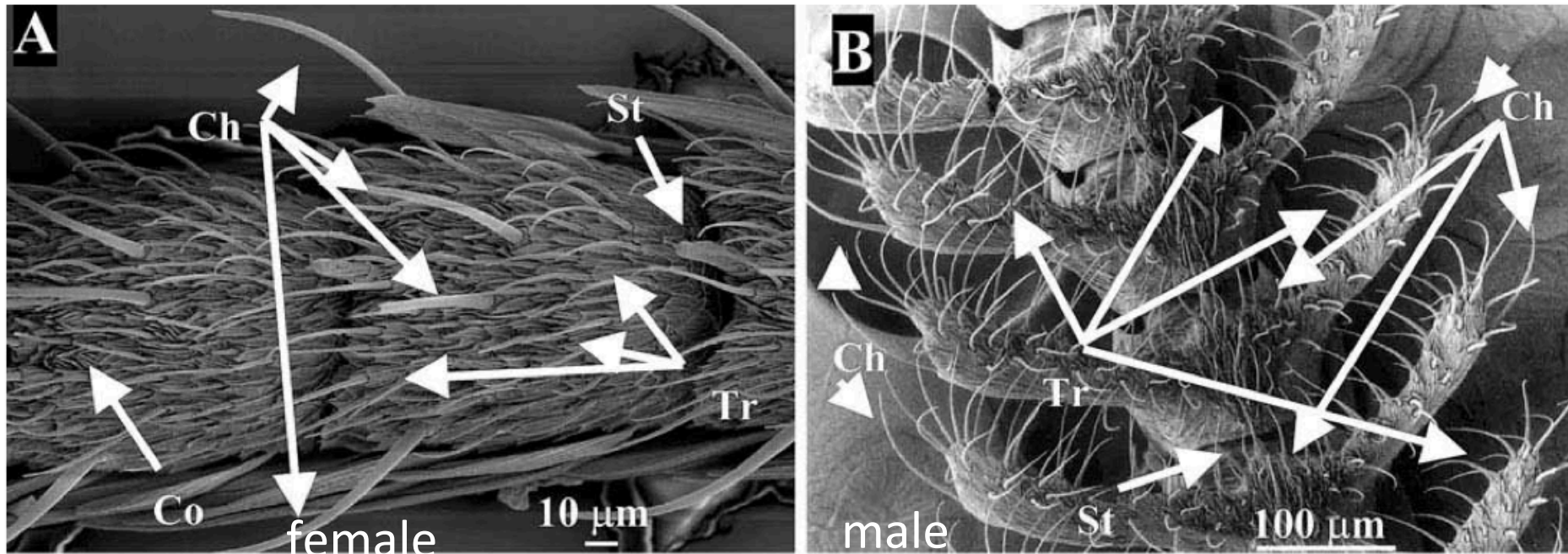
integrating research and capacity building

- Animals reared at icipe and sequencing at ILRI
- Transcriptomes from 10 different tissues
- Assembly and annotation jamboree in Nairobi 2 weeks ago



What Have We Found So Far?

Genes upregulated male versus female antennae,
endogenous viruses and cases of horizontal transfer



Members of the *Busseola Fusca* Genomics Consortium

Absolomon Kihara
Alan Orth
Alistair Darby
Appolinaire Djikeng
Bruno Leru
Caitlin Miller ('14)

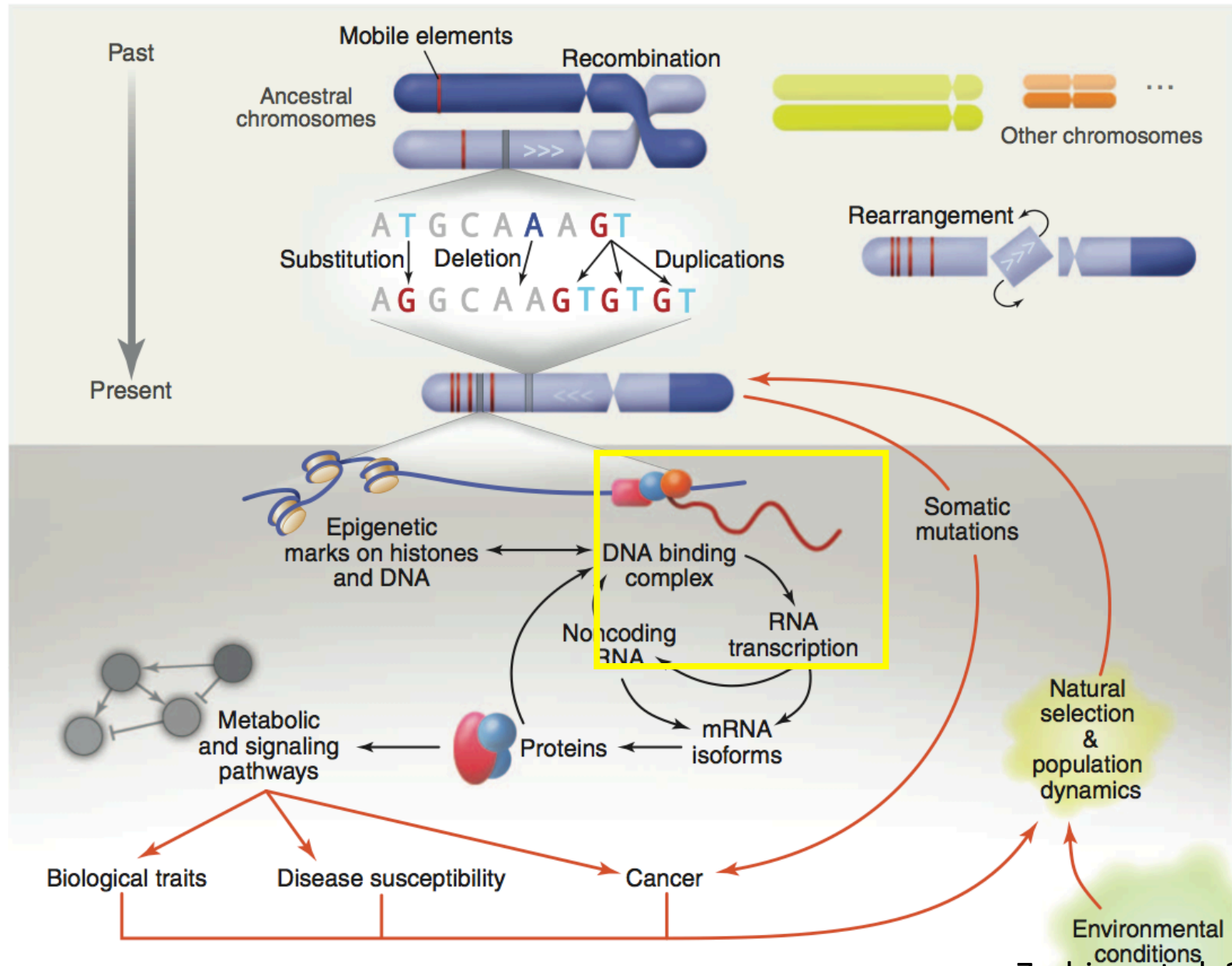
Okalang Uthman
Wacoo Paul
Alexander Ssamula
Muzoora Saphan
Kerfua Susan
Tolbert Sonda

Elizabeth Ajema Luvai
Ernest Lutomia
Chebon Lorna Jemosop
Ngalah Stephen Bidii
Damaris Matoke
Gachara Grace

Harrison Ndung'u
Wachiuri Kelvin
Mwangi
Obange Faith
Towett Sharon
Chepkemoi



Questions in Biology Range from Simple to Complex



Many Studies Focus on Understanding Differences in Gene Expression

Within a Species

Different cells with different jobs

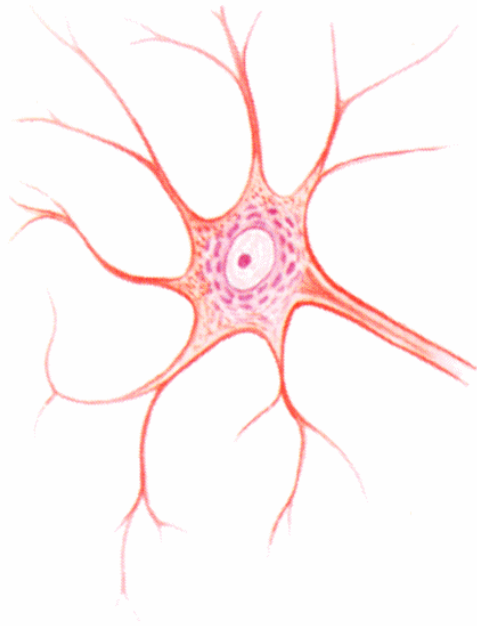
- Neurons
- Skin cells
- Intestinal cells
- Liver cells
- Blood cells
- Fat cells
- Germ cells
- Muscle cells

Same
DNA!

Responses to different conditions

- Heat/Cold
- Food/Hunger
- Disease
- Stress
- Danger
- Development
- Aging
- Damaging chemicals

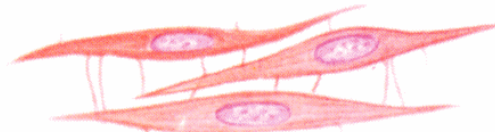
NERVE CELL



MUSCLE CELLS



Striated (voluntary)



Smooth (involuntary)

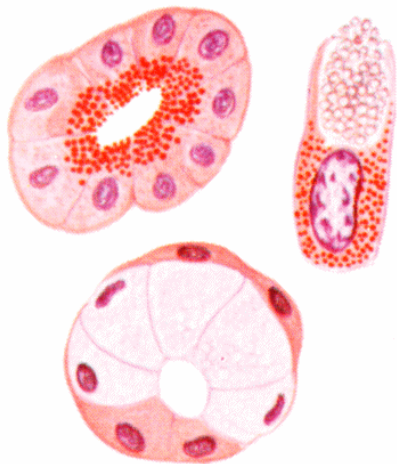


Cardiac

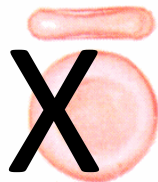
BONE CELL



GLAND CELLS

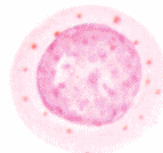


BLOOD CELLS

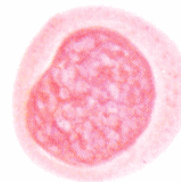


Red blood cells

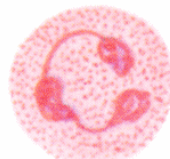
White blood cells



Lymphocyte



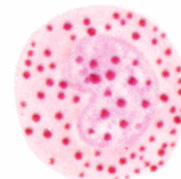
Monocyte



Neutrophil



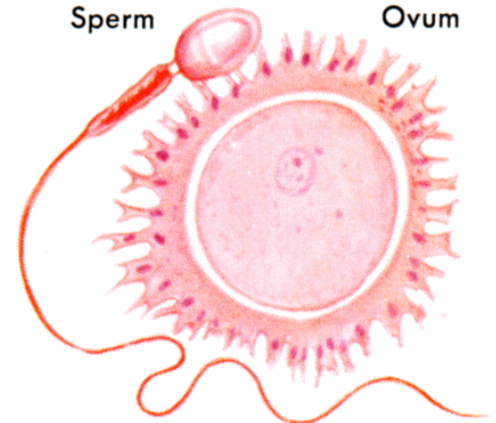
Eosinophil



Basophil

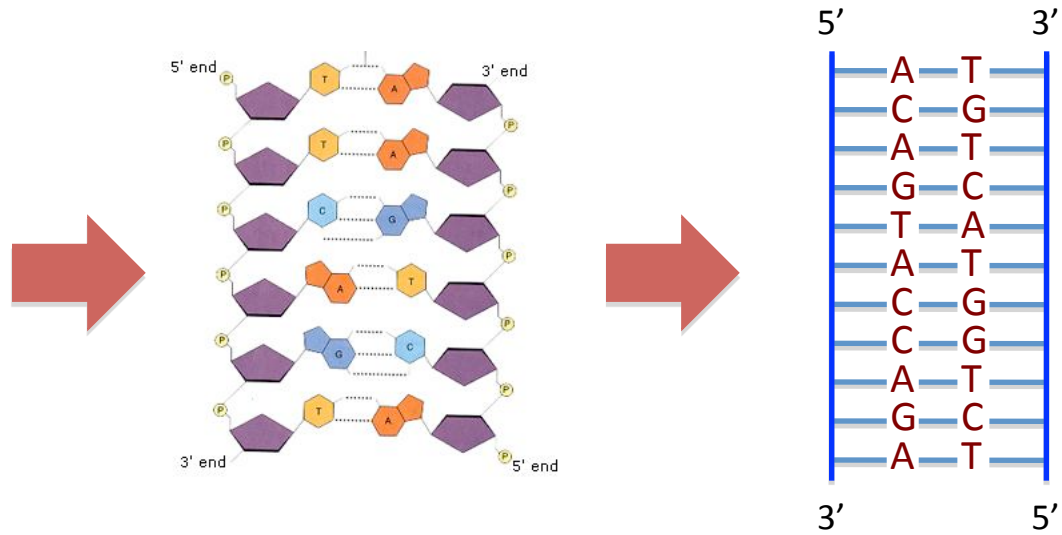
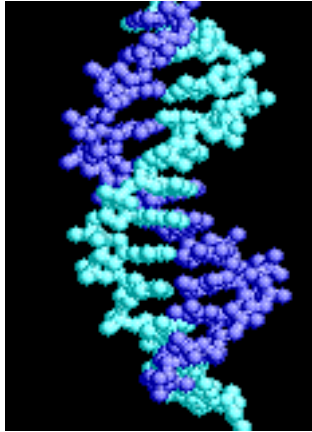
Sperm

Ovum



REPRODUCTIVE CELLS

DNA

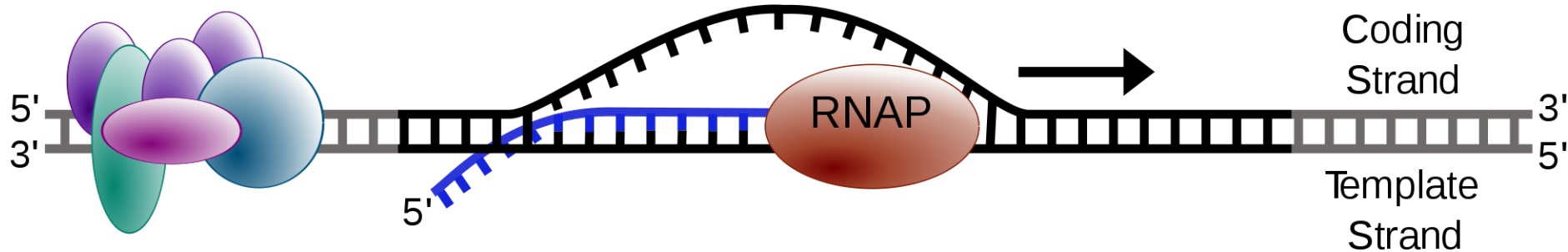


5' ACAGTACCAGACCAGACCATAACATACCATC 3'
3' TGTCATGGTCTGGTCTGGTATGTATGGTAG 5'

ACAGTACCAGACCAGACCATAACATACCATC

RNA

Transcription produces mRNA polymers

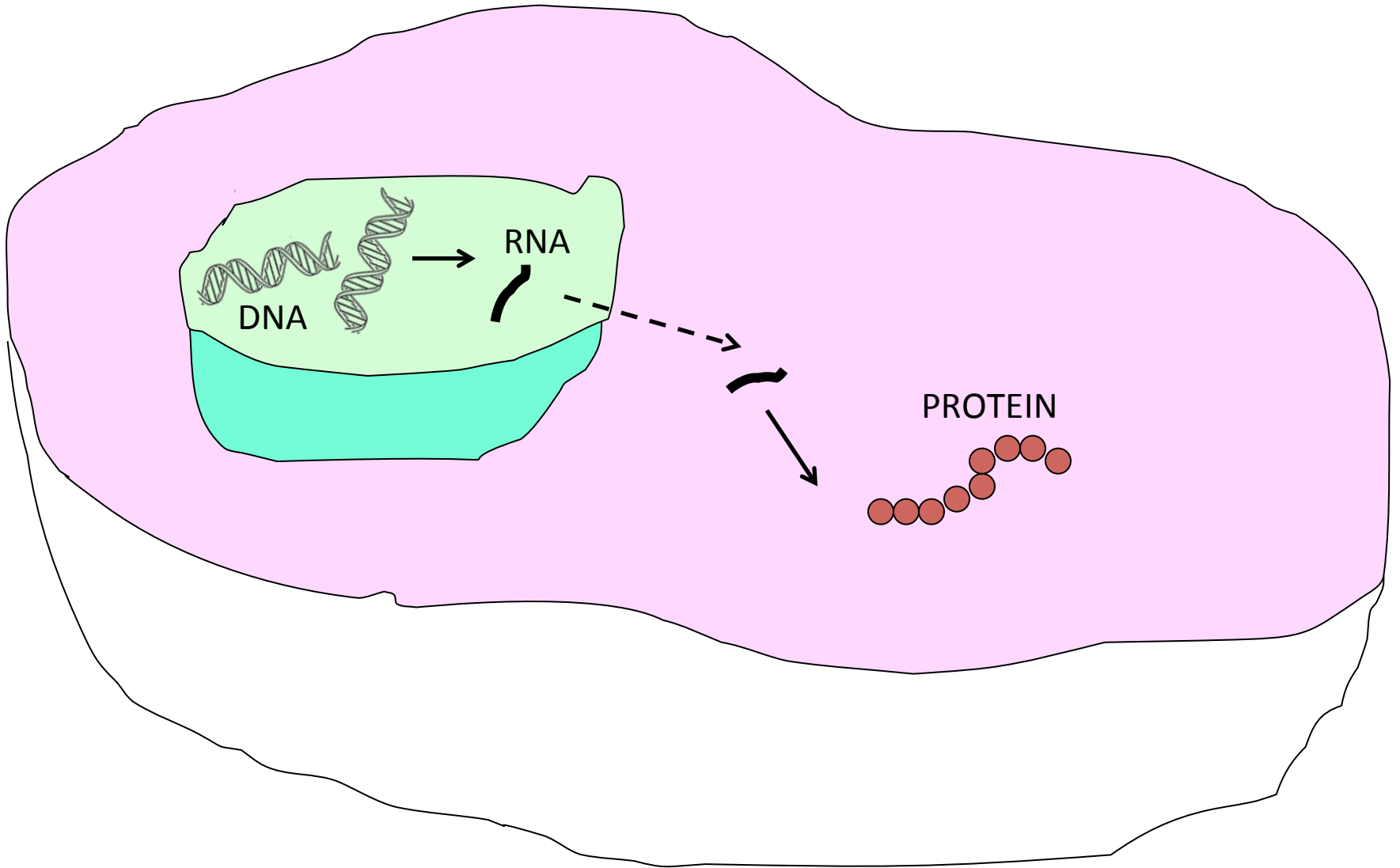


Also rRNAs, tRNA, siRNA, etc.

RNA is typically single-stranded, but can fold back on itself, making small hairpins, and therefore is capable of forming complex structures

Thus, mRNAs are messengers, but many other RNA molecules are functional and act much like proteins do

THE CENTRAL DOGMA



Coding DNA → proteins → us

The genetic code

		Second letter					
		U	C	A	G		
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA Stop UAG Stop	UGU } Cys UGC } UGA Stop UGG Trp	U C A G	Third letter
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G	
	A	AUU } Ile AUC } AUA } AUG Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G	
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }	U C A G	

Synonymous vs. nonsynonymous differences

Proline
four-fold degenerate amino acid

C C U
C C C
C C A
C C G

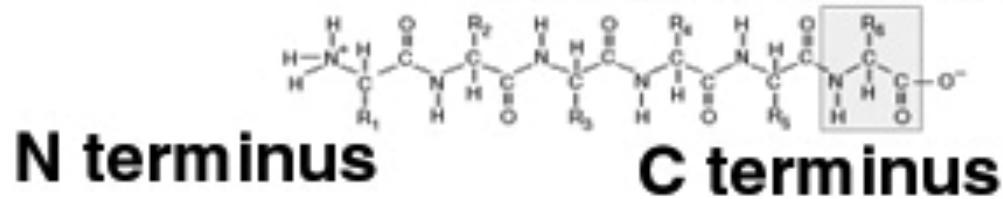
 Synonymous changes
 Nonsynonymous changes

Arginine

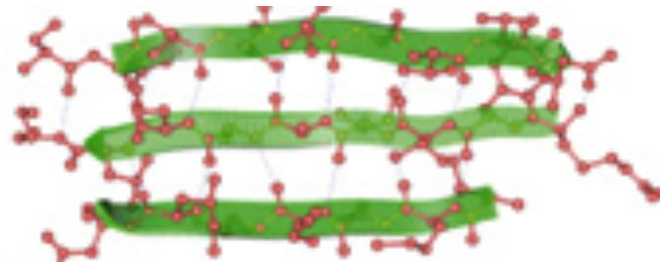
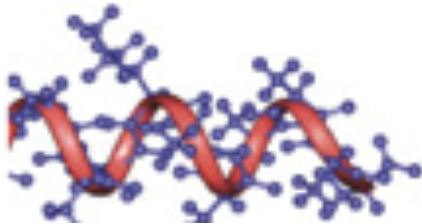
C  T

Protein Structure

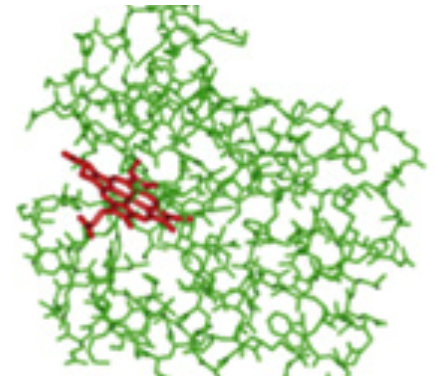
1^o structure



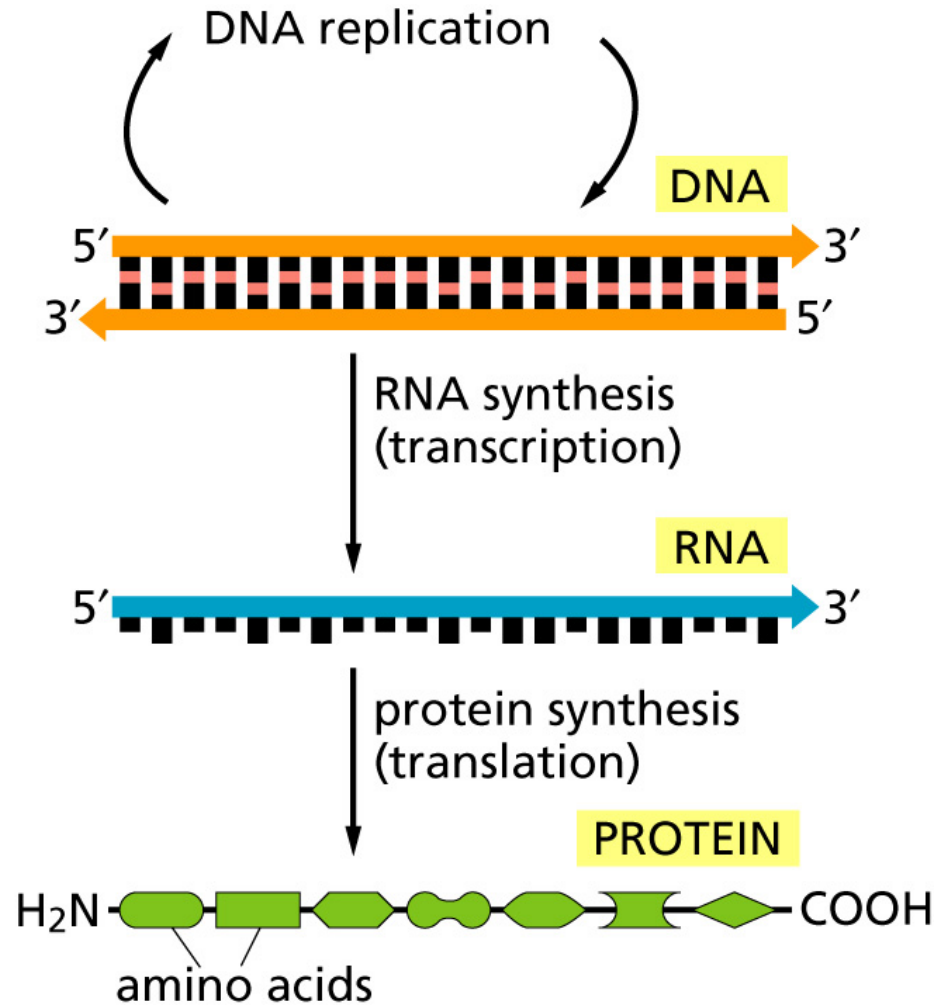
2^o structure



3^o structure

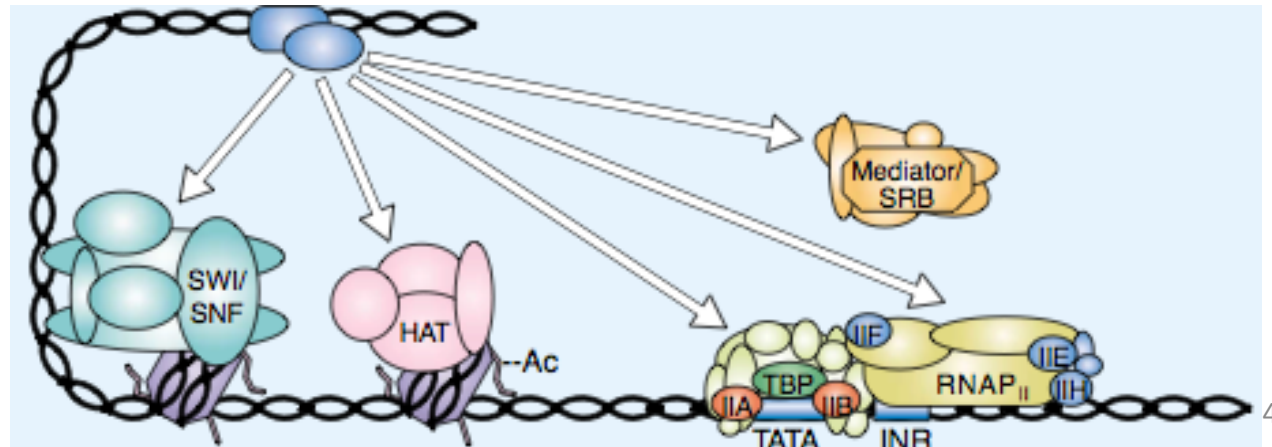
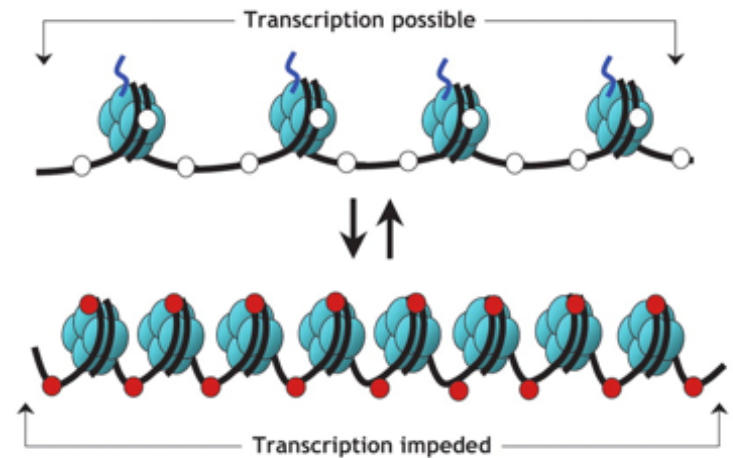


Central Dogma

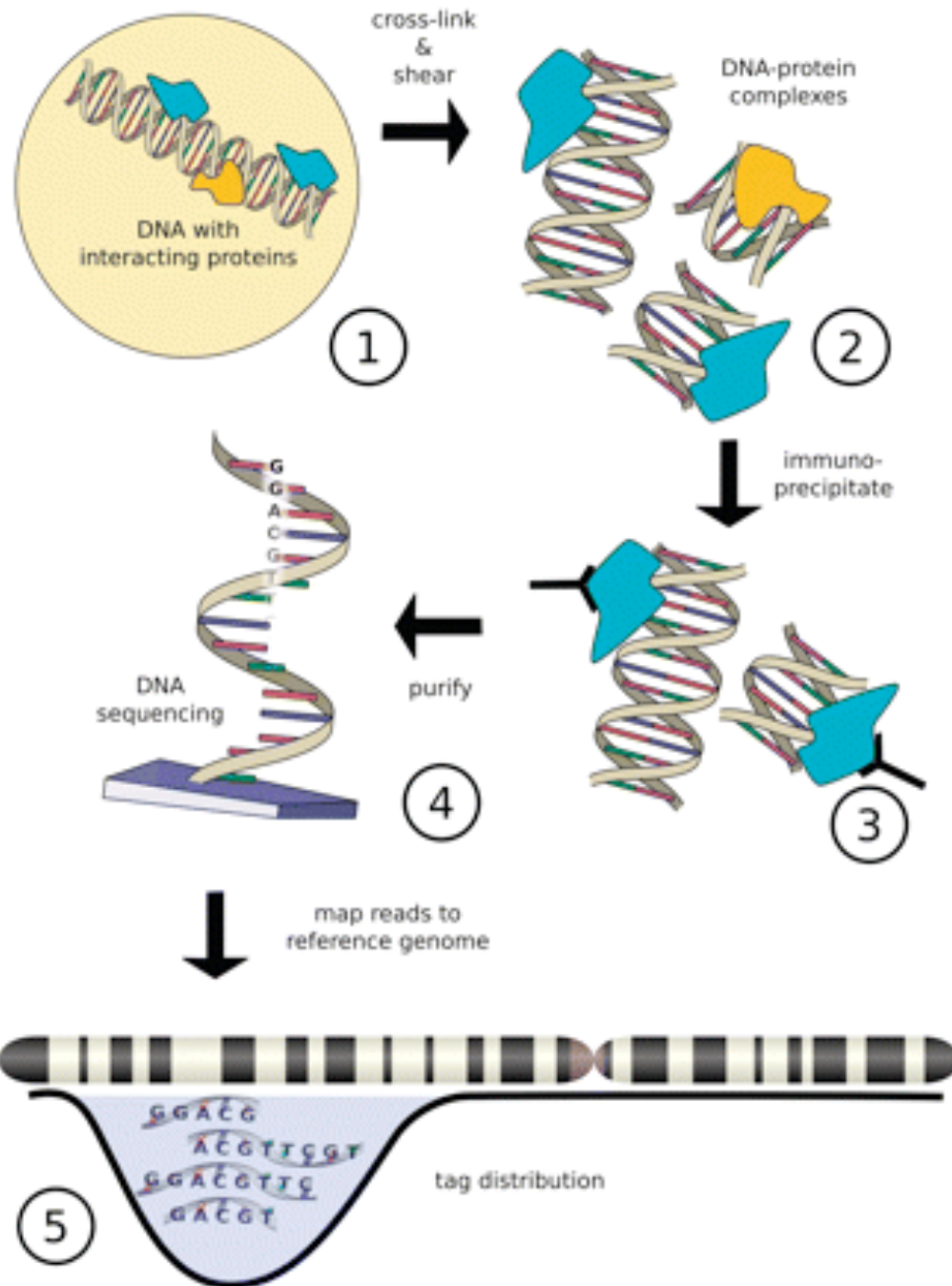


How Can Transcription Levels Be Quantified Using Bioinformatics and Genomics?

- DNA availability (methylation and acetylation)
 - Bisulfite sequencing
- Transcription Initiation and Elongation
 - CHIP-Seq
- Transcript termination, processing and export
 - RNA-Seq



ChIP-Seq



- Take cells of interest and cross-link (fix proteins to the DNA that they bind to while “in action”)
- Shear DNA
- Immunoprecipitate YPOI (your protein of interest, here the blue protein) and discard everything else
- Release the protein from the fragments it was bound to
- Sequence the DNA
- Map the data to the genome to see what genes nearby might be regulated by YPOI

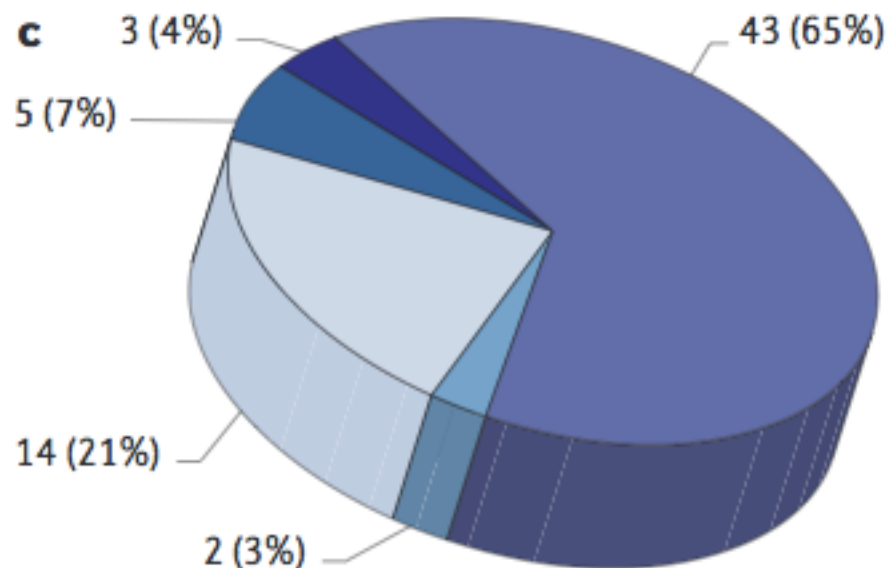
Believe it or not: how much can we rely on published data on potential drug targets?

Florian Prinz, Thomas Schlange and Khusru Asadullal

A recent report by Arrowsmith noted that the success rates for new development projects in Phase II trials have fallen from 28% to 18% in recent years, with insufficient efficacy being the most frequent reason for failure (Phase II failures: 2008–2010. *Nature Rev. Drug Discov.* 10, 328–329 (2011))¹. This indicates the limi-

to 'feasible/mar of pursuing a development project could ultimately cost billions of Euros. Even in activities such as drug discovery programmes

results that are published are hard to reproduce. However, there is an imbalance between this apparently widespread impression and its public recognition (for example, see REFS 2,3), and the surprisingly few scientific publications dealing with this topic. Indeed, to our knowledge, so far there has been no published



- Inconsistencies
- Not applicable
- Literature data are in line with in-house data
- Main data set was reproducible
- Some results were reproducible

Editorial

Ten Simple Rules for Reproducible Computational Research

Geir Kjetil Sandve^{1,2*}, Anton Nekrutenko³, James Taylor⁴, Eivind Hovig^{1,5,6}

1 Department of Informatics, University of Oslo, Blindern, Oslo, Norway, **2** Centre for Cancer Biomedicine, University of Oslo, Blindern, Oslo, Norway, **3** Department of Biochemistry and Molecular Biology and The Huck Institutes for the Life Sciences, Penn State University, University Park, Pennsylvania, United States of America, **4** Department of Biology and Department of Mathematics and Computer Science, Emory University, Atlanta, Georgia, United States of America, **5** Department of Tumor Biology, Institute for Cancer Research, The Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, Norway, **6** Institute for Medical Informatics, The Norwegian Radium Hospital, Oslo University Hospital, Montebello, Oslo, Norway

Replication is the cornerstone of a cumulative science [1]. However, new We further note that reproducibility is just as much about the habits that ensure than to do it while underway). We believe that the rewards of reproducibility will

10 Simple Rules for Reproducibility

- Rule 1: For Every Result, Keep Track of How It Was Produced
- Rule 2: Avoid Manual Data Manipulation Steps
- Rule 3: Archive the Exact Versions of All External Programs Used
- Rule 4: Version Control All Custom Scripts
- Rule 5: Record All Intermediate Results, When Possible in Standardized Formats
- Rule 6: For Analyses That Include Randomness, Note Underlying Random Seeds
- Rule 7: Always Store Raw Data behind Plots
- Rule 8: Generate Hierarchical Analysis Output, Allowing Layers of Increasing Detail to Be Inspected
- Rule 9: Connect Textual Statements to Underlying Results
- Rule 10: Provide Public Access to Scripts, Runs, and Results