

Computational approaches to study Genetics

遺伝学を研究するための計算アプローチ

2020/03/09

RIKEN

Dr. Jeffrey Fawcett



Arithmer Seminar

弊社社員によるスライドではなく、弊社にて平均週一回で開催されている、「Arithmer Seminar」にて登壇頂いた外部の方々によるスライドです。

Computational approaches to study Genetics

Jeffrey Fawcett

(31 Oct 2019 @Arithmer)





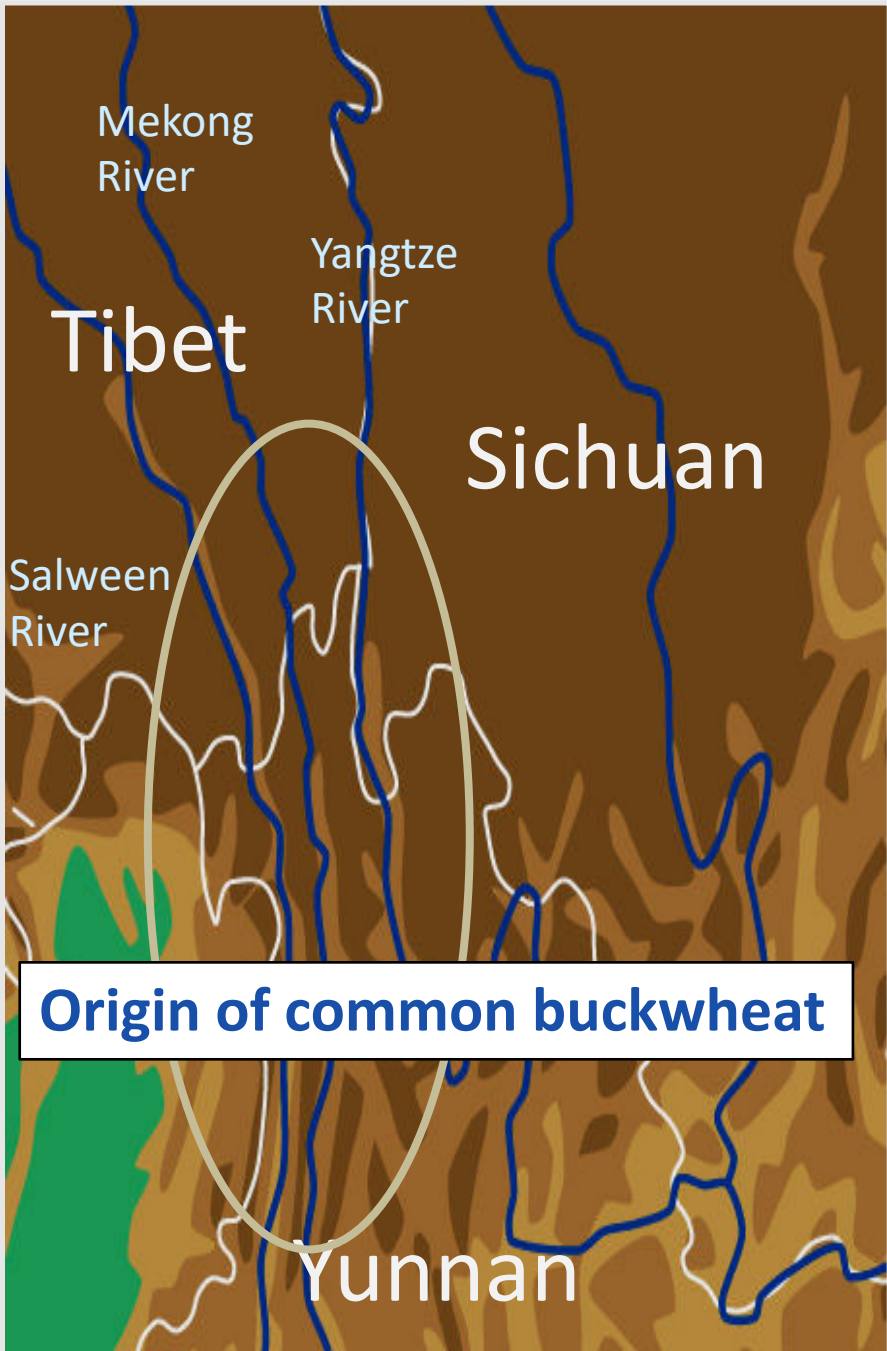
梅里雪山 (6,740m)
Mt. Meili Xueshan



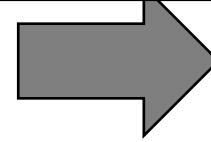
CHINA







Domestication



Fagopyrum esculentum
ssp. *ancestrale*

Fagopyrum esculentum
ssp. *esculentum*

**Yasuo Yasui
(Kyoto Uni)**



**Chengyun Li
(Yunnan Agricultural Uni)**



**Jeffrey Fawcett
(RIKEN iTHEMS)**



**Takanori Ohsako
(Kyoto Pref Uni)**



**Diane Lister
(Cambridge Uni)**

Aim of Buckwheat Project

- ❑ Collect & maintain genetic resources of buckwheat-related species
(many wild species are facing extinction due to development)
- ❑ Understand domestication process of buckwheat
- ❑ Identify genes important for breeding of buckwheat

“Extensive characterization of domestication-related genes in buckwheat by utilizing the genetic resource of Yunnan province, China”

KAKENHI Fostering Joint International Research B, PI: Yasuo Yasui (Kyoto U)

中国雲南省の野生ソバ遺伝資源を活用した栽培化関連遺伝子の網羅的同定

科研費・国際共同研究強化B 研究代表者：安井康夫（京都大）

Why am I involved?



lots of DNA data to be generated



A
T
C
G

owlication.com

```
ATGGGCCAAGTTTTTGAAGTCTTAA  
ATTTAAAAATCATATACACGTTGTA  
AAAATGCGTAGGTTTCATGAATGAAA  
TCAGAATTTACCACACTTACTGAAA  
ATGACTTGTGAGTTGTGATGGATAA  
GTTTGATTAAATAAGAAAGGTAAT  
GCATATGGCTACAAAATGAAAGATT
```

- Huge advance in technology to generate large-scale DNA data
- Need for people with computational skills and knowledge of genetics (very few such people that work on buckwheat, horses, etc)

Main Research Interest

- Process of genetics/evolution responsible for creating the diversity of life
- Apply existing knowledge to (e.g. agronomically important) species

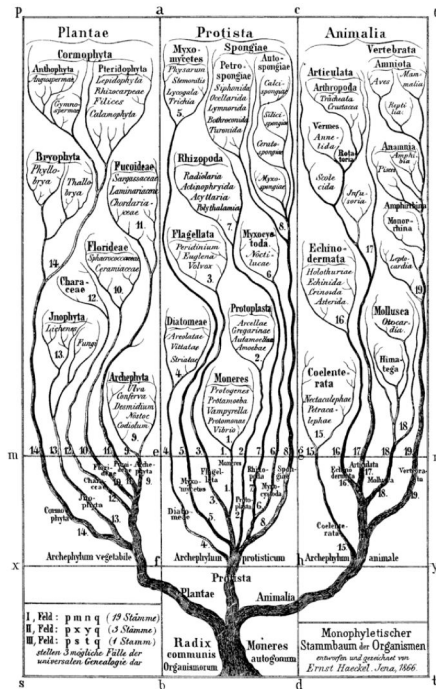
Outline

- Basic concepts of genetics and evolution
- What is written in the DNA?
- How can we know what's written in the DNA?
- How can we associate “genotype” with “phenotype”?
 - Research on Thoroughbred horses

History of Abstraction in Biology

Biology is... complex, diverse, changes over time

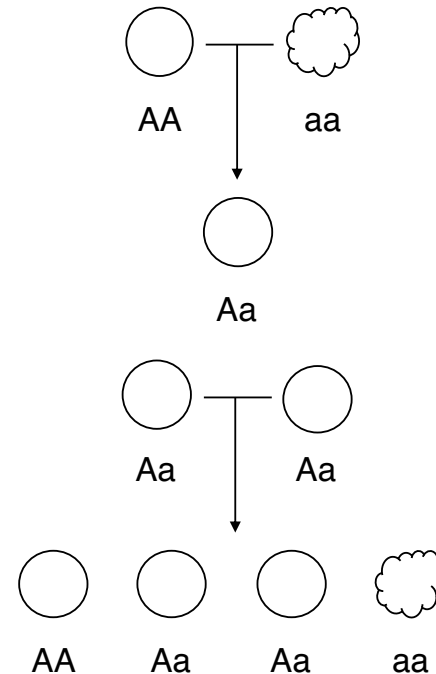
Evolution: “Tree of Life”



Branching process

“Difference” increases with time

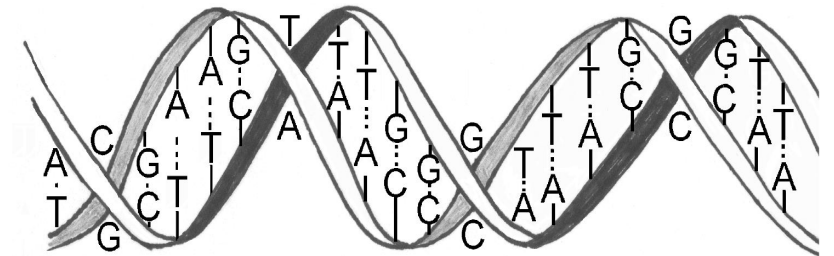
Mendelian Genetics



Traits as symbols (AA, Aa, aa)

Change of frequency as a probabilistic process

DNA

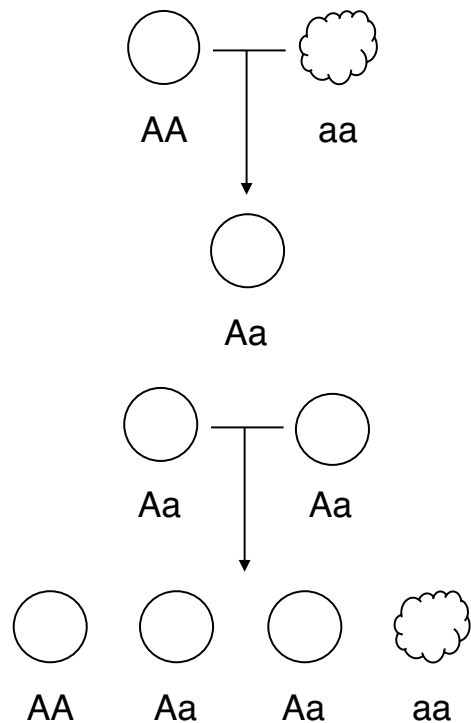


Forms of life as digital information, “sequence” of letters

Basic concepts of genetics

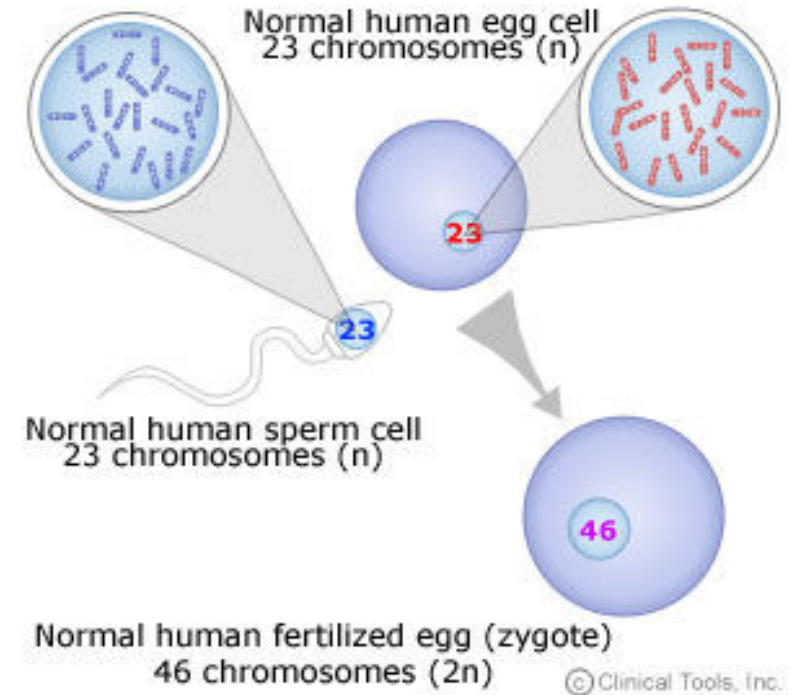
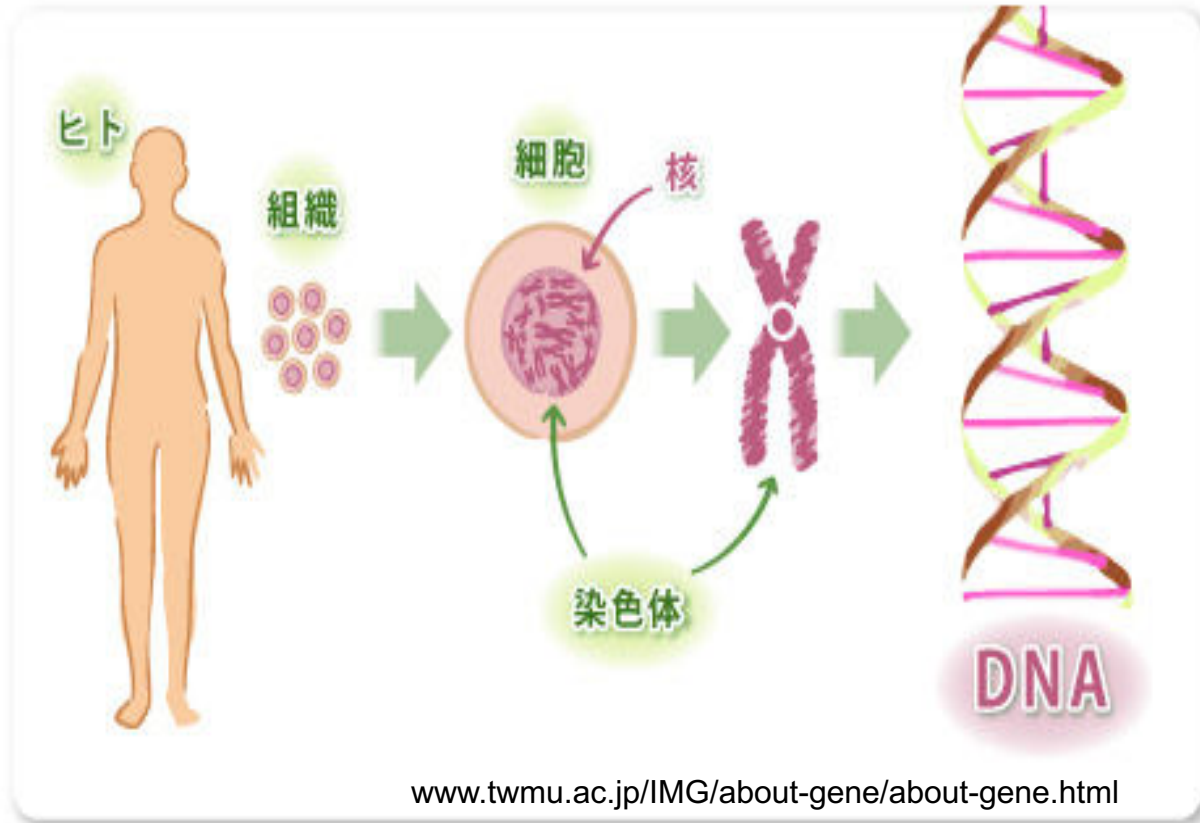
- Parent and offspring are similar – genetic information is passed on

What is passed on and how?



- Information of A and a is transmitted across generations without changing
- The “phenotype” (appearance) of AA and Aa are the same

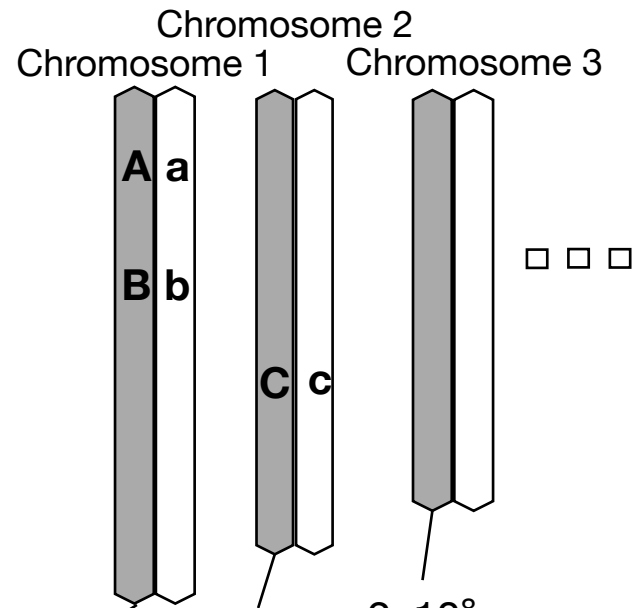
Mechanism of Heredity



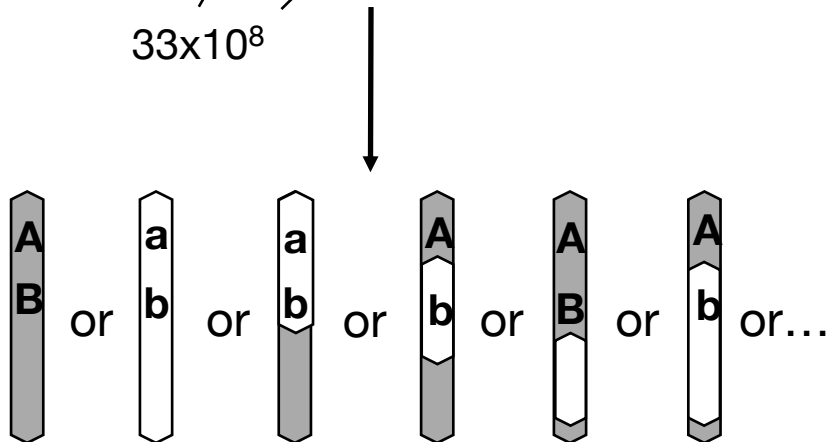
Genome: entire set of genetic information

Every human has 2 genomes (1 from each parent)

Mechanism of Heredity

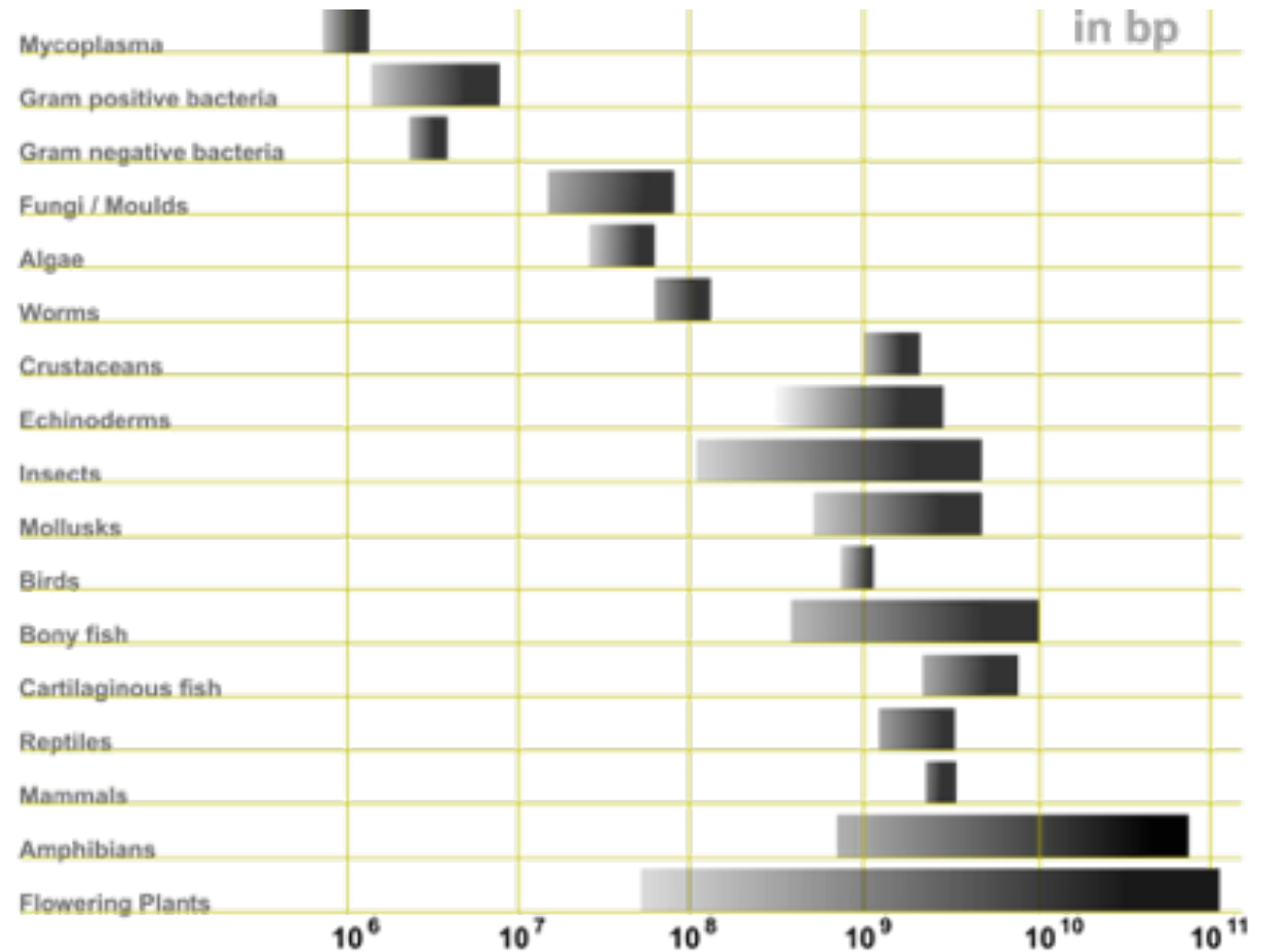


human: 23x2 chromosomes (22x2 + XX or XY),
 $\sim 33 \times 10^8$ (3 billion) (x2) nucleotides



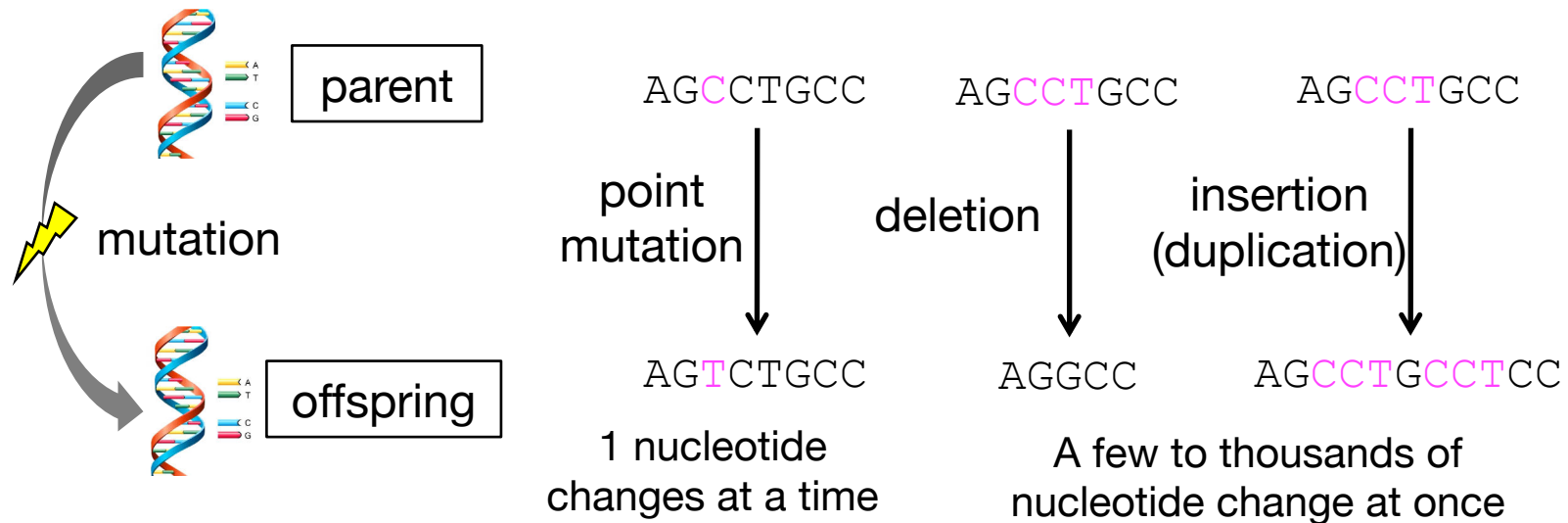
Large variation in genome size

Species	#Nucleotides (x10 ⁸)
Budding yeast	0.121
Fruit fly	1.75
Fugu	3.9
Rice	4.5
Maize	25
Mice	27
Human	33
Onion	150
Grasshopper	650
Lungfish	1300
Canopy plant (キヌガサソウ)	1490



Mutation

“Replication” cannot generate diversity, genetic information changes

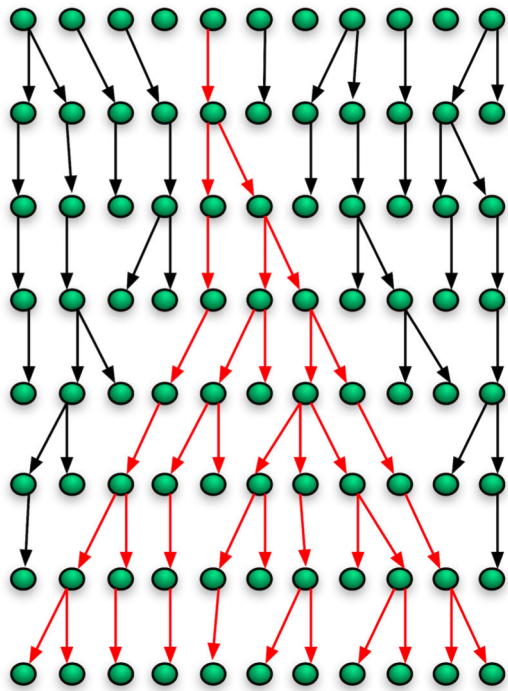


~100 mutations per generation in humans

~1/10⁸ point mutations per nucleotide per generation

Evolution

- Heredity
- Mutation => some of them change the “phenotype”
- “Population process” (competition, “struggle for survival”)



Many individuals do not produce any offspring

We can't observe mutations in those individuals

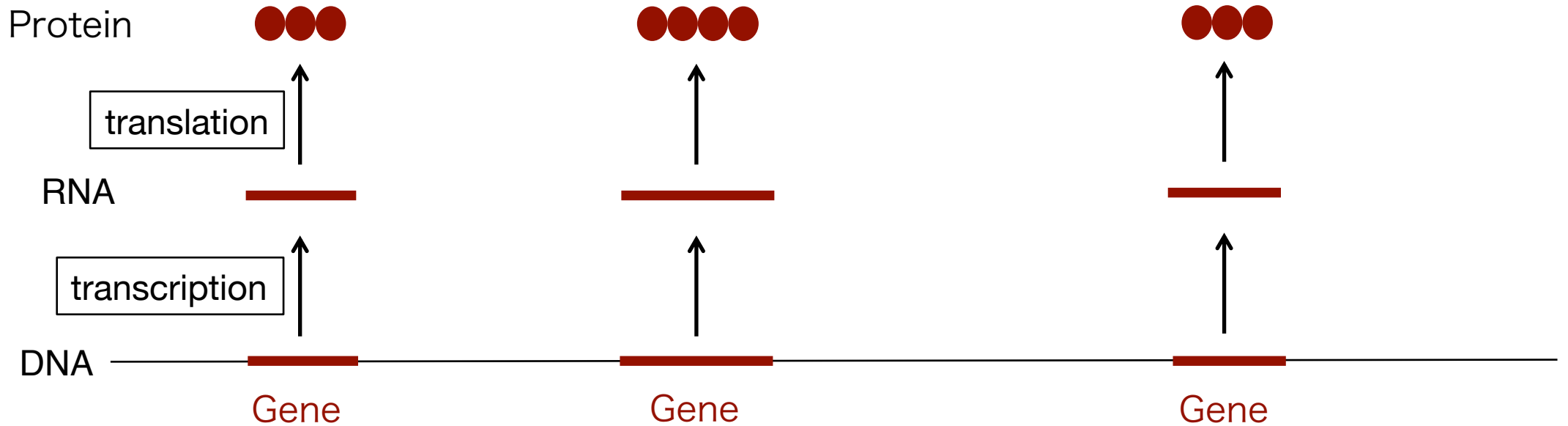
Do all individuals have equal chance to produce offspring?

Did the mutation change the chance to produce offspring?

Outline

- ❑ Basic concepts of genetics and evolution
- ❑ What is written in the DNA?
- ❑ How can we know what's written in the DNA?
- ❑ How can we associate “genotype” with “phenotype”?
 - Research on Thoroughbred horses

What is written in the genome?



All cells contain the same set of genes

The amount and timing of RNA/proteins produced differ
(very complicated process)

RNA -> Amino Acid (translation)

- 20 kinds of amino acids

DNA -> RNA
(transcription)

A (Adenine) -> A (Adenine)

G (Guanine) -> G (Guanine)

C (Cytosine) -> C (Cytosine)

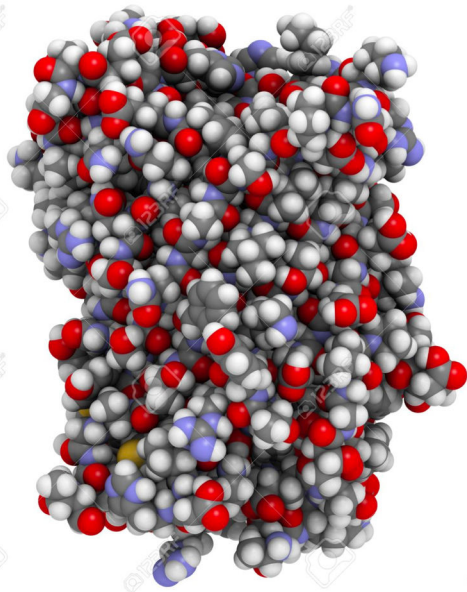
T (Thymine) -> U (Uracil)

		Second Letter																			
		U		C		A		G													
1st letter	U	UUU Phe	UCU Ser	UAU Tyr	UGU Cys	U	UUC Leu	UCC Ser	UAC Stop	UGC Stop	C	UUA Leu	UCA Stop	UAA Stop	UGA Stop	A	UUG Trp	UCG Trp	UAG Stop	UGG Trp	G
	C	CUU Leu	CCU Pro	CAU His	CGU Arg	U	CUC Leu	CCC Pro	CAC Gln	CGC Arg	C	CUA Leu	CCA Pro	CAA Gln	CGA Arg	A	CUG Leu	CCG Pro	CAG Gln	CGG Arg	G
	A	AUU Ile	ACU Thr	AAU Asn	AGU Ser	U	AUC Ile	ACC Thr	AAC Lys	AGC Arg	C	AUA Met	ACA Thr	AAA Lys	AGA Arg	A	AUG Met	ACG Thr	AAG Lys	AGG Arg	G
	G	GUU Val	GCU Ala	GAU Asp	GGU Gly	U	GUC Val	GCC Ala	GAC Glu	GGC Gly	C	GUA Val	GCA Ala	GAA Glu	GGA Gly	A	GUG Val	GCG Ala	GAG Glu	GGG Gly	G
						3rd letter															

redundancy in genetic code

Proteins as amino acid sequences

Protein



www.123rf.com

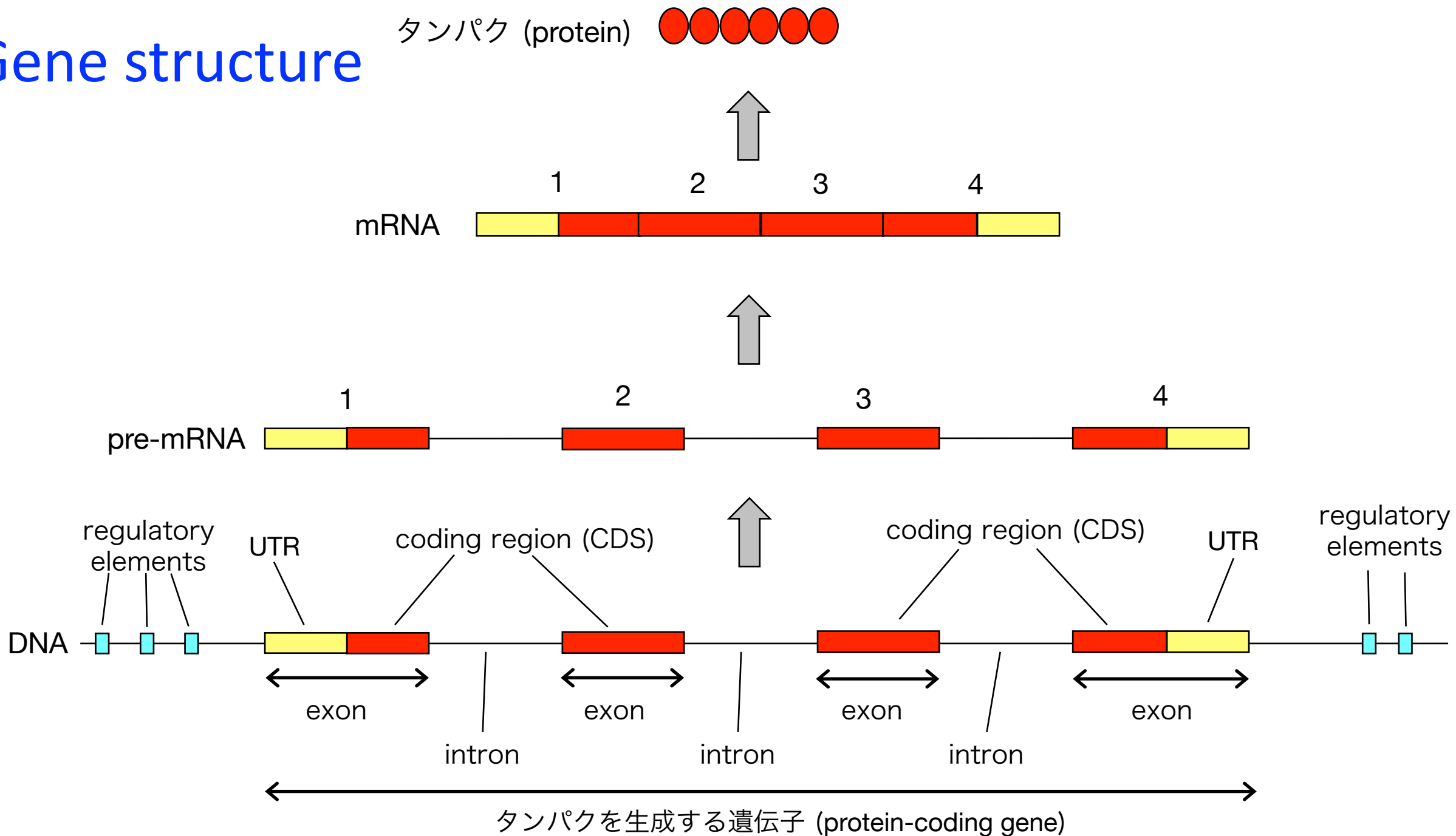
Amino Acid sequence of Histone H4 gene

Human	YEE TR G V L K V F L E N V I R D A V T Y T E H A K R R K T V T A M D V V Y A L
Fruitfly	YEE TR G V L K V F L E N V I R D A V T Y T E H A K R R K T V T A L D V V Y A L
Tomato	YEE TR G V L K I F L E N V I R D S V T Y T E H A R R K T V T A M D V V Y A L
Volvox	YEE TR T V L K N F L E N V I R D S V T Y T E H A R R K T V T A M D V V Y A L
Budding Yeast	YEE V R A V L K S F L E S V I R D S V T Y T E H A K R R K T V T S L D V V Y A L

- The “same” genes are similar across species
- Similar sequence -> similar function

➤ Many genes share high similarity across distantly related species

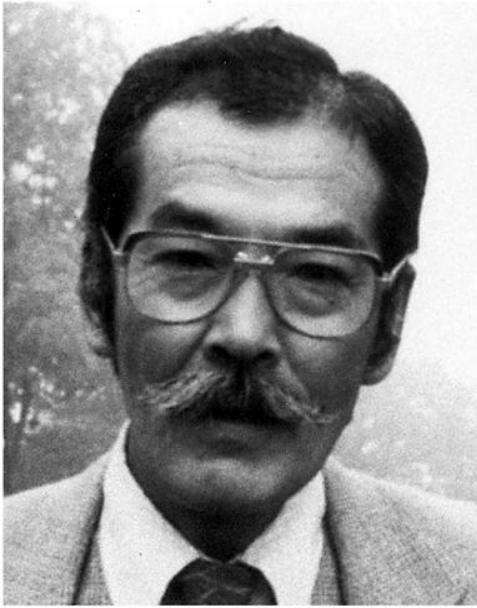
Gene structure



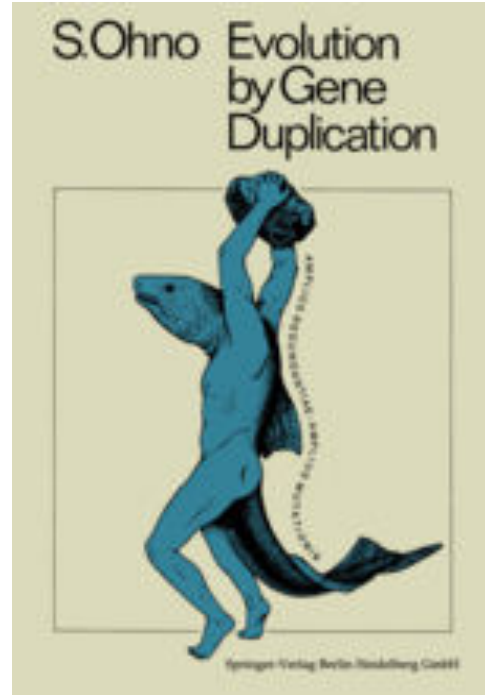
Number of genes are not so different across species

Species	#Nucleotides (x10 ⁸)	#Genes
Budding yeast	0.121	6,000
Fruit fly	1.75	17,000
Fugu	3.9	28,000
Rice	4.5	40,000
Maize	25	32,000
Mice	27	23,000
Human	33	21,000
Onion	150	?
Grasshopper	650	?
Lungfish	1300	?
Canopy plant (キノガサソウ)	1490	?

Evolution by Gene Duplication



Susumu Ohno



- most genes were created by duplication of another existing gene
- mutation can create a new function while keeping the original function

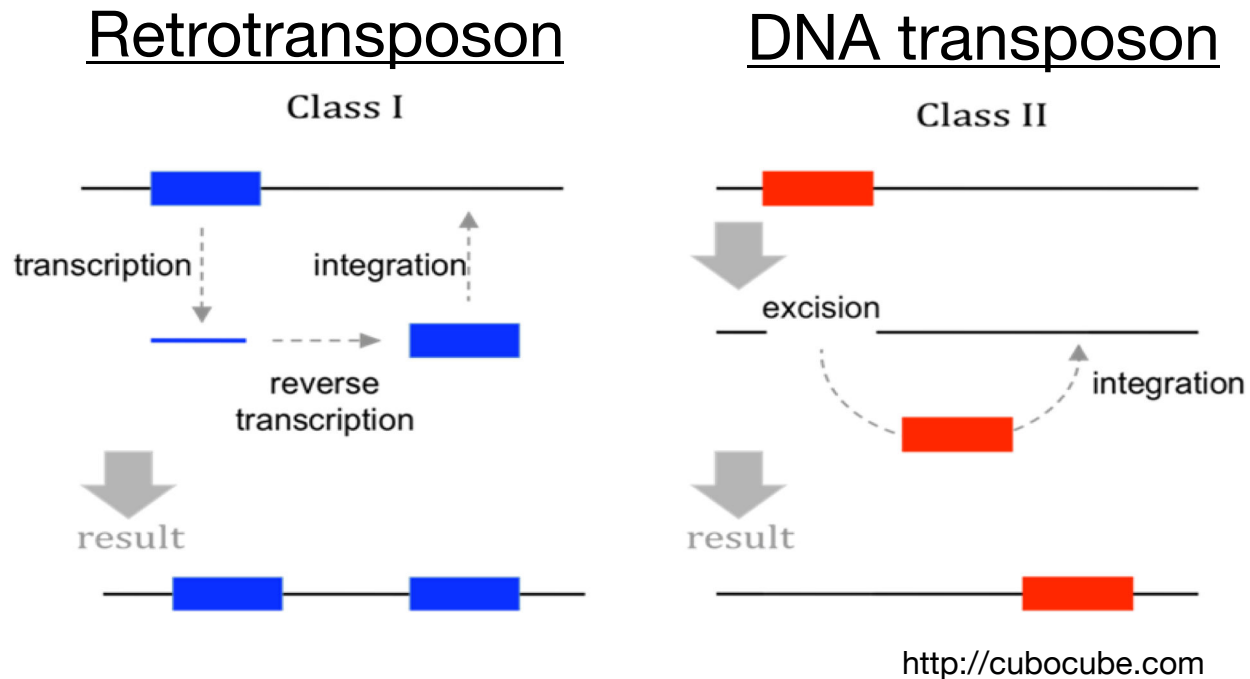
What is written in the genome?

□ Human genome

protein-coding sequences: 1-2%

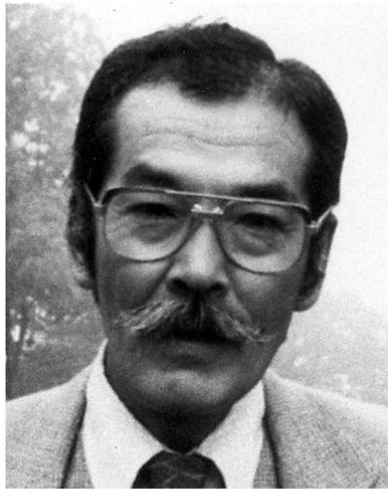
regulatory sequences: ~5-10%

Transposable Elements: >60%



- major reason for genome size variation
- virus-like parasites that amplify
- have their own “genes”
- often harmful, sometimes beneficial

A lot of unnecessary information in the genome!?



Susumu Ohno

“So much *Junk DNA* in our Genome”
(Susumu Ohno, 1972)

We all acquire ~100 new mutations but most of us are fine (i.e. most mutations are harmless)

“Junk DNA” should be removed during evolution!?!?

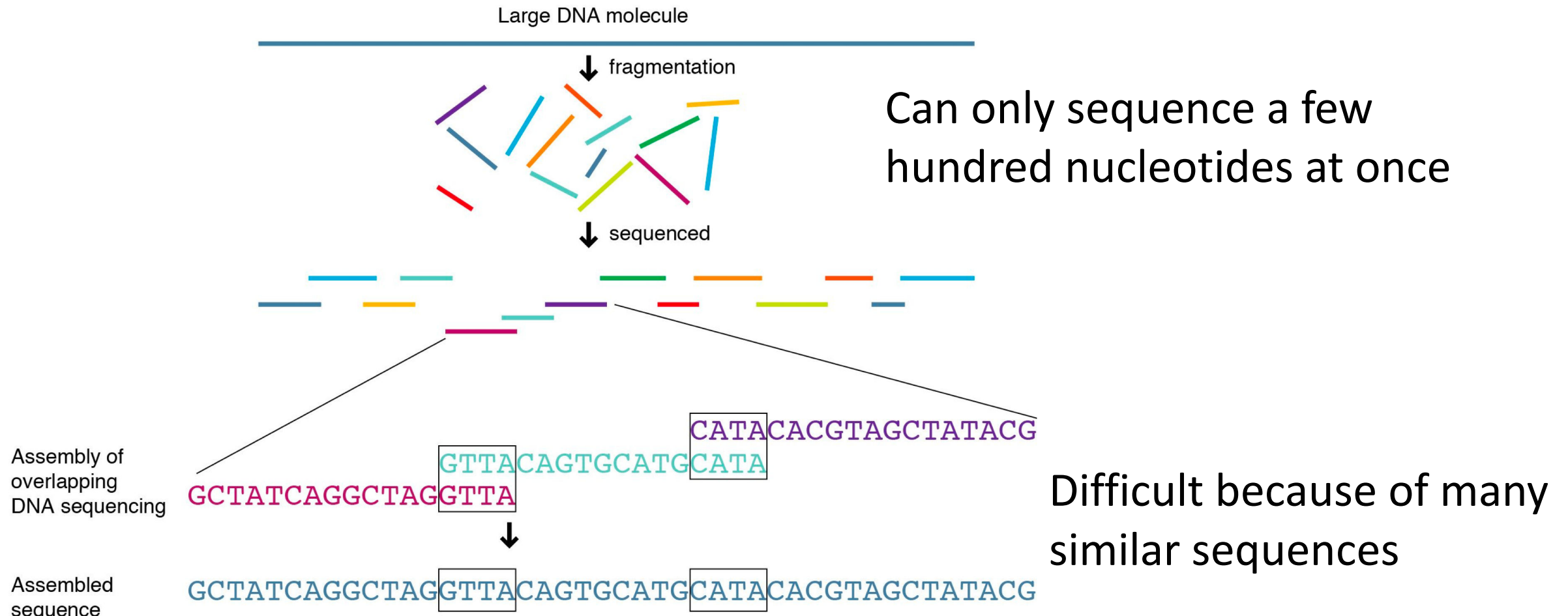
- each “junk” is removed but lots of “junk” are being generated
- they are parasites that try to survive themselves

Outline

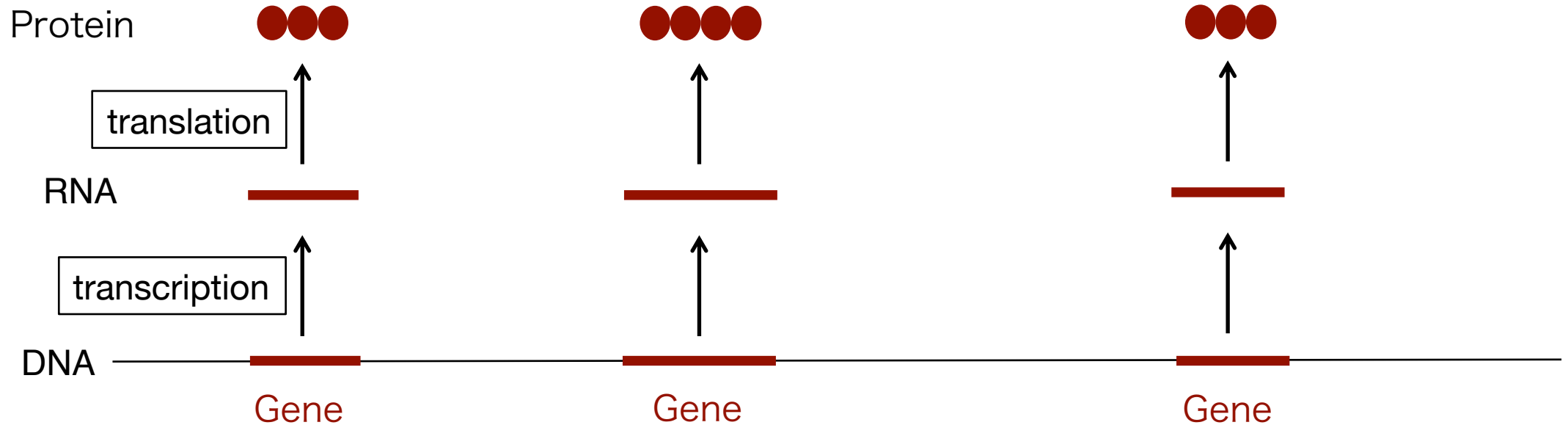
- ❑ Basic concepts of genetics and evolution
- ❑ What is written in the DNA?
- ❑ How can we know what's written in the DNA?
- ❑ How can we associate “genotype” with “phenotype”?
 - Research on Thoroughbred horses

How do we know the sequence of the genome?

- Genome sequencing and assembly



How can we “find” the genes?



Extract RNA, sequence, and “map” them to the genome

Not so simple because of many similar genes and the presence of introns

How can we “find” the genes?

- Histone H4 gene

Human	MSGRGKGGKGLGKGGAKRHRK V LRDNIQGITKPAIRRLARRGGVKRISGLIYEE T R G VL K V F LE N VIRD A VTYTEHA K RKTVT A MDVVYALKRQGRTLYGF G
Fruitfly	M T GRGKGGKGLGKGGAKRHRK V LRDNIQGITKPAIRRLARRGGVKRISGLIYEE T R G VL K V F LE N VIRD A VTYTEHA K RKTVT A L D VVYALKRQGRTLYGF G
Tomato	M S GRGKGGKGLGKGGAKRHRK V LRDNIQGITKPAIRRLARRGGVKRISGLIYEE T R G VL K I FLE N VIRD S VTYTEHA R RKTVT A MDVVYALKRQGRTLYGF G
Volvox	M S GRGKGGKGLGKGGAKRHRK V LRDNIQGITKPAIRRLARRGGVKRISGLIYEE T R T VL K N FLE N VIRD S VTYTEHA R RKTVT A MDVVYALKRQGRTLYGF G
Budding Yeast	M S GRGKGGKGLGKGGAKRHRK I LRDNIQGITKPAIRRLARRGGVKRISGLIYEE V R A VL K S FLE S VIRD S VTYTEHA K RKTVT S L D VVYALKRQGRTLYGF G

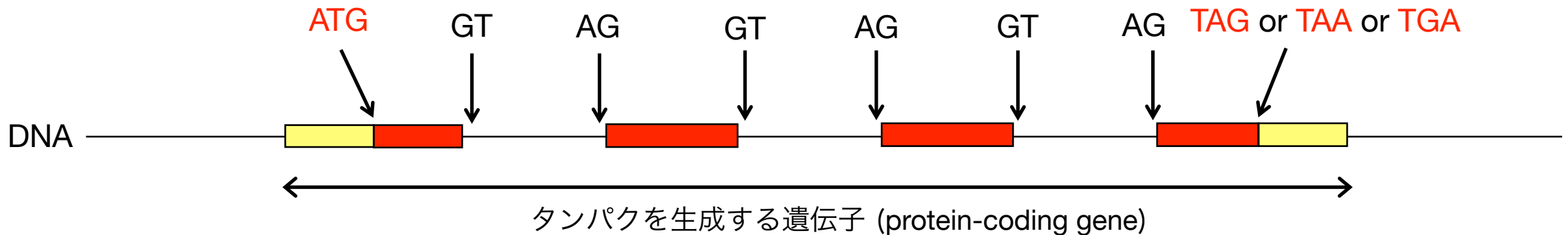
Search for similar sequences in the genome

Most genes in 1 species have similar genes in other species

Similar genes are likely to have similar functions

How can we “find” the genes?

Predict genes based on their features



Nucleotide composition is statistically different

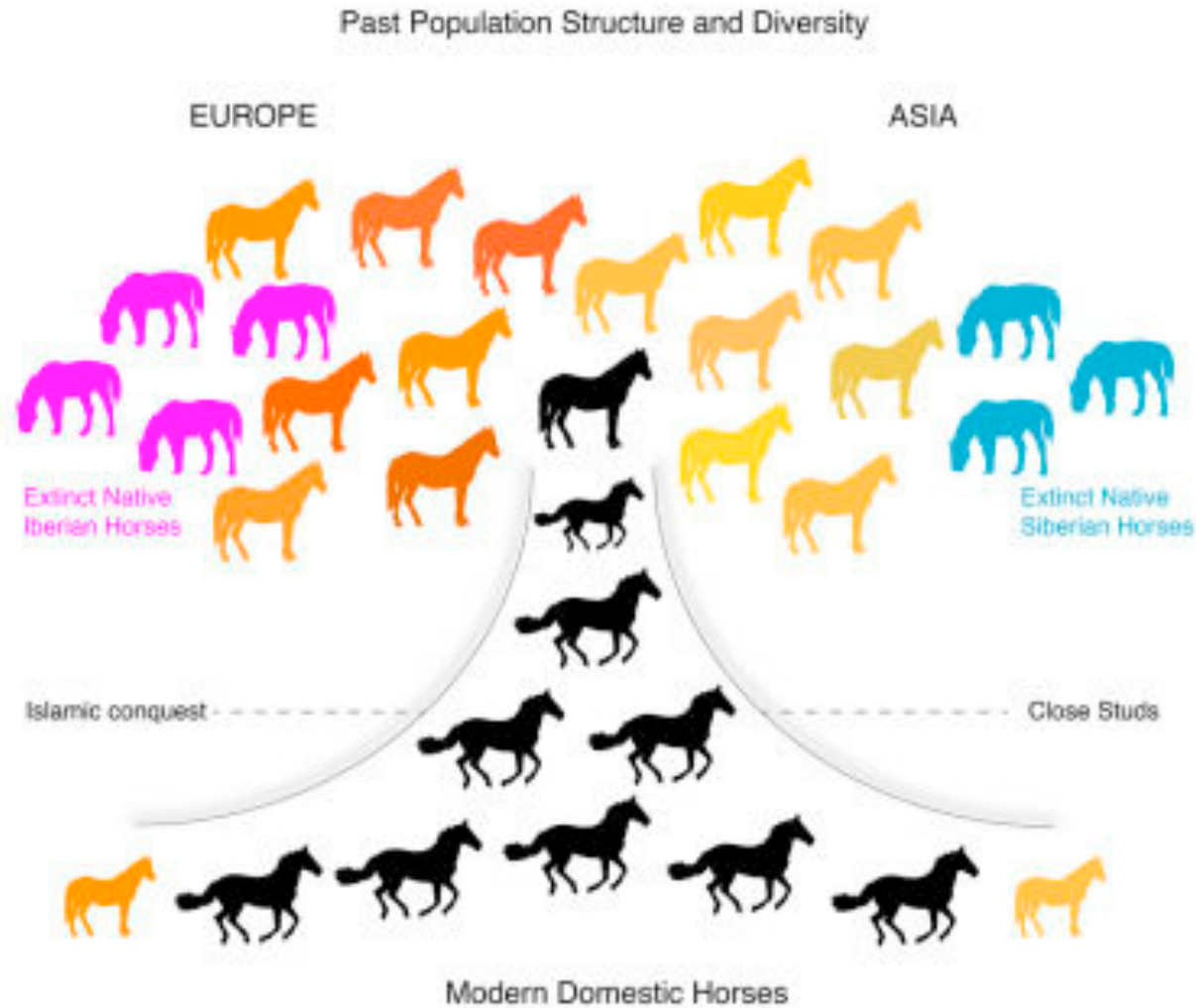
protein-coding sequence vs non-coding sequences, introns

GT-AG of intron vs GT-AG not associated with introns

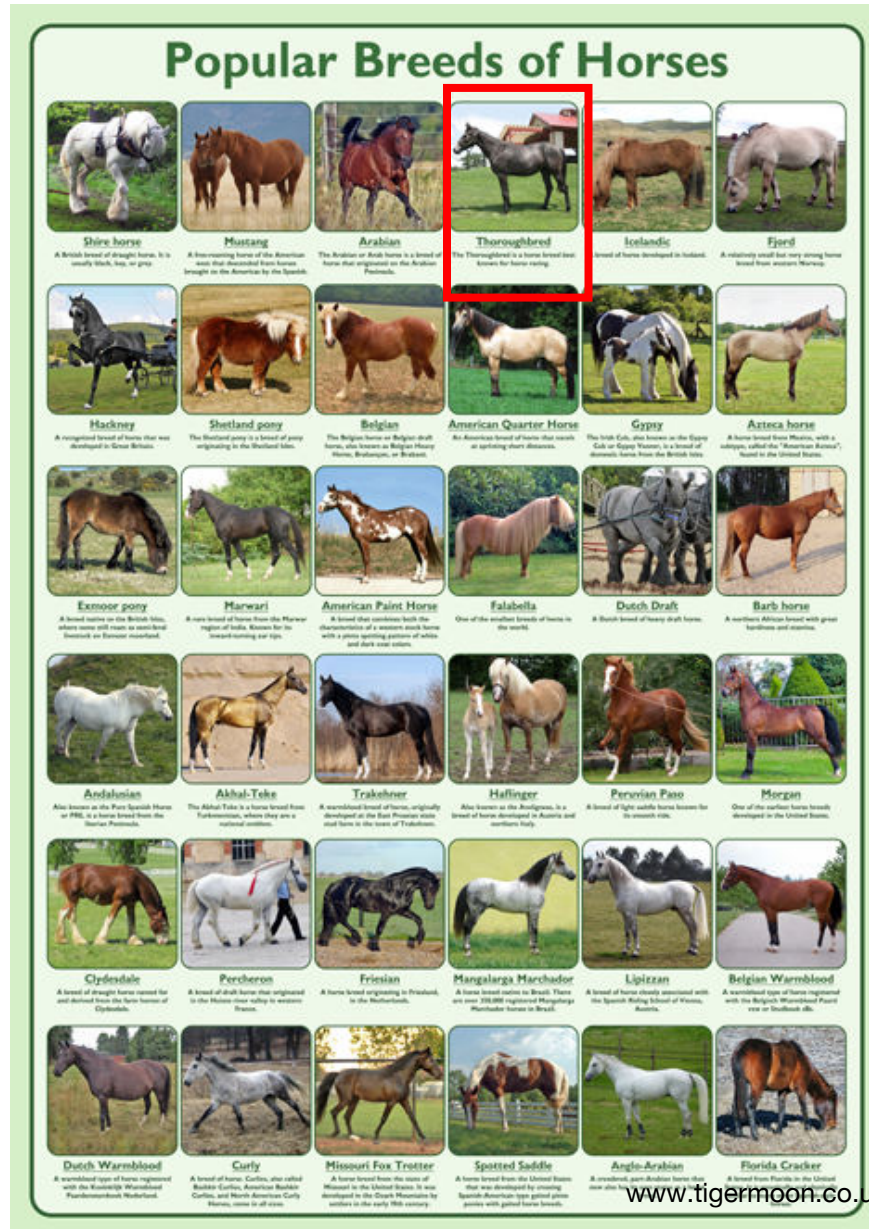
Outline

- ❑ Basic concepts of genetics and evolution
- ❑ What is written in the DNA?
- ❑ How can we know what's written in the DNA?
- ❑ How can we associate “genotype” with “phenotype”?
 - Research on Thoroughbred horses

Domestication of Horses (~5500 yrs ago)



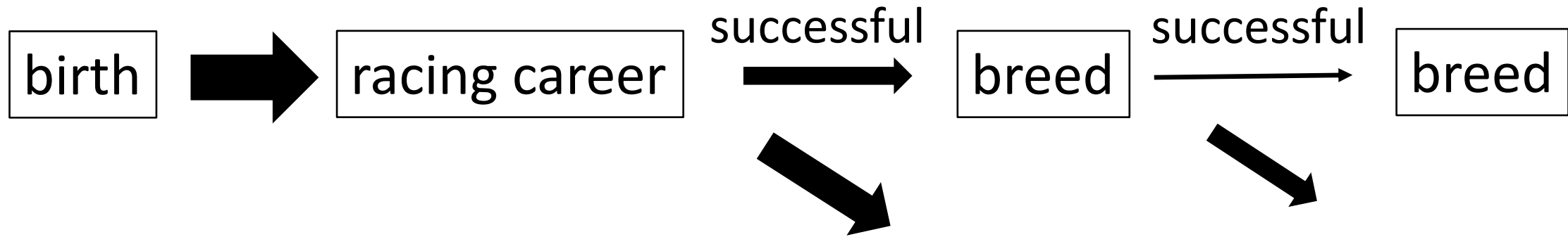
Many diverse breeds established after domestication



Thoroughbreds

- Originated in 18th Century (3 “founder” stallions)
- Horses that win races are selected to breed in many different countries

Selective breeding of Thoroughbreds in Japan



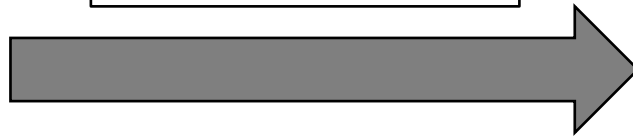
- ~7000 horses are born and registered at JRA (Japan Racing Association)
- Only 10-20 males per generation (>50% of females) are selected for breeding
- Thoroughbreds are much faster than other horse breeds due to selective breeding for 20-30 generations
=> but, genetic information has not been utilized yet

Many differences between wild and cultivated Buckwheat

Fagopyrum esculentum
ssp. ancestrale



Domestication

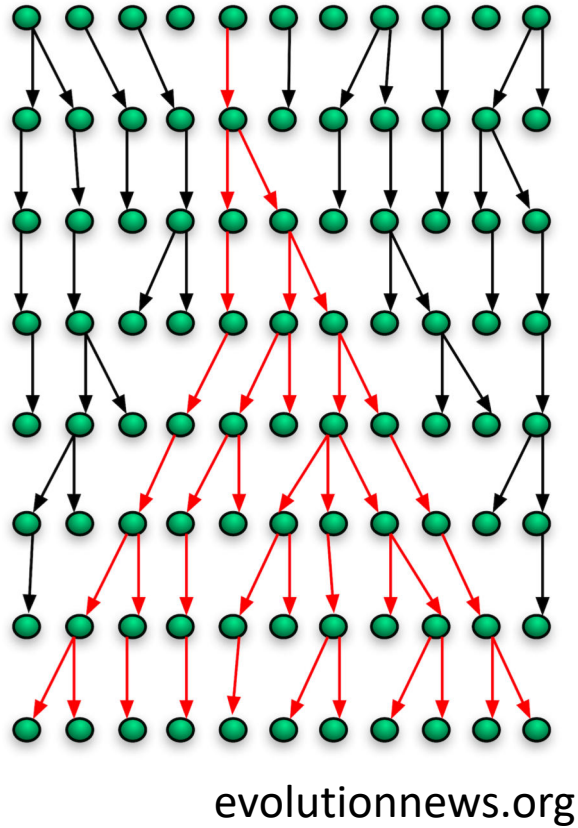


Loss of seed shattering
Loss of seed dormancy
Erect growth
Larger seeds

Fagopyrum esculentum
ssp. esculentum



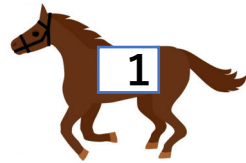
Process of domestication



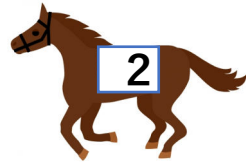
- Mutations occur that result in a desirable trait (might already exist in wild, or occur afterwards)
- Humans preferentially select and breed those individuals
- **Can we identify the mutations selected by humans?** (should speed up the selective breeding)
- **Can we identify other useful mutations/genes that could be useful for breeding?**

Variation in the DNA of each individual

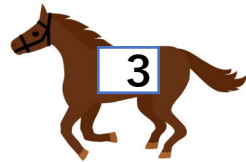
Single Nucleotide Polymorphism (SNP)



ATGTACG **T**GAGTCAG **T**ACGATGCAGAT **C**TATGAAGA **C**ATCCGA...
ATGTACG **T**GAGTCAG **T**ACGATGCAGAT **G**TATGAAGA **A**ATCCGA...



ATGTACG **T**GAGTCAG **G**ACGATGCAGAT **G**TATGAAGA **C**ATCCGA...
ATGTACG **C**GAGTCAG **G**ACGATGCAGAT **G**TATGAAGA **C**ATCCGA...

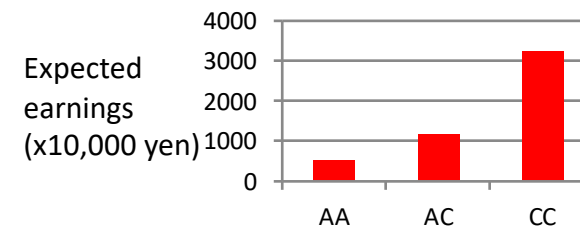


ATGTACG **C**GAGTCAG **T**ACGATGCAGAT **G**TATGAAGA **A**ATCCGA...
ATGTACG **C**GAGTCAG **G**ACGATGCAGAT **C**TATGAAGA **A**ATCCGA...

SNP of Myostatin gene

TT: Long distance
CT: Middle distance
CC: Short distance

↑ SNP of Gene X related to racing ability



Example of variation (SNP) associated with racing ability

第1号掲載の企画が馬主・生産者の間でも大反響

ミオスタチン遺伝子(スピード遺伝子)最新研究報告企画・第2弾

遺伝子検査は競馬を変える!?

~めざせ「テーラーメイド調教」~

C/C(短距離型)、C/T(中距離型)、T/T(中-長距離型)という3パターンのミオスタチン遺伝子型により、馬の距離適性傾向が予測できる。小誌第1号ではその研究報告について特集したところ、競馬サークル内外から多数の反響を頂戴した。そこで今回も、日本におけるこの研究の第一人者である戸崎晃明博士に、この分野の最新研究状況などについてご寄稿をいただいた。

文／公益財団法人 競走馬理化学研究所 遺伝子分析室 戸崎晃明

1969年生まれ、栃木県出身。昭和大学大学院博士課程を修了。昭和大学博士(薬学)、京都大学博士(農学)、昭和大学医学部兼任講師、内閣古農家大学特任教授、岐阜大学応用生物科学部委員獣医学系教授、公益財団法人競走馬理化学研究所に入所後は、研究部に籍を置き、遺伝地図作製やゲノム解読などの国際プロジェクト「Horse Genome Project」に参加する。最近問題になりつつある「遺伝子ドーピング」の問題にも取り組んでいる。




図1 ミオスタチン遺伝子は筋肉量を抑制する遺伝子

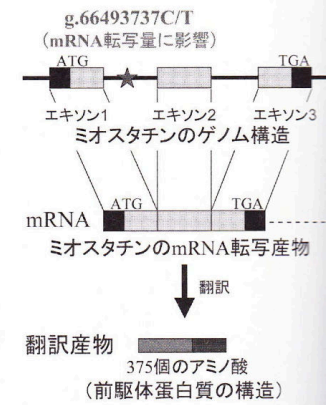
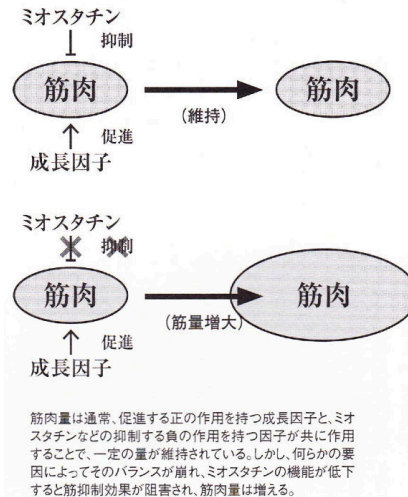


図2 g.66493737C/T多型とミオスタチン遺伝子の構造

資料B ミオスタチン遺伝子型別の距離適性傾向

ミオスタチン	筋量傾向	戸崎博士統計	UCD統計
C/C型	やや多い	1000~1800m	1000~1600m
C/T型	普通	1200~2000m	1400~2400m
T/T型	やや少ない	1800m以上	2000m以上

SNP in myostatin gene affects optimum racing distance

Example of variation (SNP) associated with racing ability

【天皇賞・秋】春連覇のフェノメノ、遺伝子検査では中距離型だった！

2014年10月28日6時0分 スポーツ報知



注目されるG1で、牡馬、牝馬の強豪古馬が

3200メートルの天皇賞・春を2連覇した。このG1でもタイトルをつかむことができるなどから、推測する材料を与えてくれるが、スピード遺伝子検査が、大きなヒントを与え

春秋合わせて区切りの150回を迎える伝統の一戦に、“春の盾”を連覇したフェノメノが、テイエムオペラオー以来史上2頭目の天皇賞3勝目を目指しエントリーした。実績から長距離向きの印象が強いが、DNAレベルから距離適性を推測する「エクイノム・スピード遺伝子検査」で、「中距離型」との判

【共同通信杯】イスラボニータ完勝“実に強い”ダービー見えた！



イスラボニータ（右）が共同通信杯を制して3連勝で皐月賞に向かう

クラシックの栄冠が見えた。降雪で代替開催（東京6日目）となった「第48回共同通信杯」が24日、東京競馬場で行われ、1番人気、蛭名正義（44）騎乗のイスラボニータが直線で抜け出して完勝。昨年のいちようS、東京スポーツ杯2歳Sに続く3連勝で東のクラシック最有力候補に浮上した。今後は皐月賞（4月20日、中山）へ直行する。【レース結果】

放牧明けで臨んだ3歳始動戦。同師は「ダービーがピークになるように調整しているが、2歳時に比べて首さしがグッとたくましくなった」と成長ぶりに目を細めた。昨年暮れには放牧先の社台ファーム（北海道千歳）から朗報が届けられた。競走馬理化学研究所が実施している遺伝子のDNA距離鑑定を受けたところ、ダービーの2400メートルまで守備範囲との評価。吉田代表は「2400メートルはやってみないと分からないが、少なくとも2000メートルまでは問題ない。シンから丈夫でキリッとした馬。これは走るぞ」。次走・皐月賞を経てダービー（6月1日、東京）へ。「さあ、これからが本当の勝負だ」。栗田博師は念願のクラシック獲りへ再び顔を紅潮させた。



1200m好走馬の次走が2200m、一見無茶な選択の裏に遺伝子検査/吉田電作マル秘週報

2014年01月29日(水)18時00分

注目数：21人

◆科学はサラブレッドの世界をどう変えていくのか

12日の京都芝外2200メートルを舞台にした3歳未勝利戦でのこと。結果は1番人気ディルガー→2番人気アグリパーバイオというごく無難な決着に終わった一方で、ある画期的な試みが実践に移されていた。

エイシンソルティエ（牝）は西園キウ舎の当世代のトップバッターを務めた馬。デビューは6月の阪神（ダ1200メートル）で0秒2差3着。2戦目（中京芝1200メートル）にはクビ差の2着まで詰め寄り、初勝利は目前と思われたが...。その後も詰めめのはなはな解消せず、前走前の時点での成績は[0-1-3-1]。とはいえ、あと一れまでと同じ短距離戦を選択するところだろう。が、西園調教1頭の舞台を選択。実に前走から1000メートルもの距離延長と

でも」ってノリでこの手の極端な条件変更をするケースもあからかに違うケース。実は“科学的な根拠”がこの1000メートル延

れているが、これが形となって表れたのが昨年の英国クラシックジャー調教師が擁するドーンアプローチだ。圧倒的なスピード入ってきたのは英ダービー。しかし研究者でもあるボルジャスピード遺伝子を強く受け継いでおり、英ダービーは適距離と究成果通りなら、愛馬は距離の壁に泣く、何ともシニカルな状況は激しく折り合いを欠き、大惨敗を喫した。

エイシンソルティエも栗東トレセンの診療所で「スピード遺伝子が向く」という判定が出た。これを受けたオーナーの指示もあけた。

Using genetic information for Thoroughbred breeding

- Can evaluate the genetic potential of each horse
=> more informed choice of which horse to keep and which to discard
- Can identify the best breeding (male x female) combination
=> more informed choice of which male and female to mate

Genome-wide SNP analysis of Japanese Thoroughbred racehorses

Jeffrey A. Fawcett^{1,2*}, Fumio Sato³, Takahiro Sakamoto¹, Watal M. Iwasaki¹,
Teruaki Tozaki⁴, Hideki Innan^{1*}



Equine Genotyping
Flexible, Accurate, Rapid Turnaround

Equine SNP70 BeadChip
Whole Genome SNP profiling

Available Exclusively From GeneSeek!

The Equine SNP genotyping BeadChip is built on Illumina's Infinium platform and is the most comprehensive solution available for genome-wide genotyping. This improved, second generation BeadChip features more than 65,000 evenly distributed SNPs derived from the EquiCabZ.0 SNP collection. The Equine BeadChip offers a powerful tool to enable identification of genes and polymorphisms that contribute to traits of interest in all major horse breeds.

Also Custom SNP Genotyping Panels Available:

- Sequenom® MassARRAY® spectrometry based detection system for sensitive, accurate, and rapid genotyping
- Easy-to-use multiplexed assay design and optimization software saves research time and helps maximize efficiency
- Flexible SNP numbers allows for economical marker assisted selection
- Parentage Testing

GENESEEEK
4665 Innovation Drive, Suite 120 • Lincoln, NE 68521
402.432.0965 • Fax: 402.432.0964
geneseeke@neogen.com • www.neogen.com

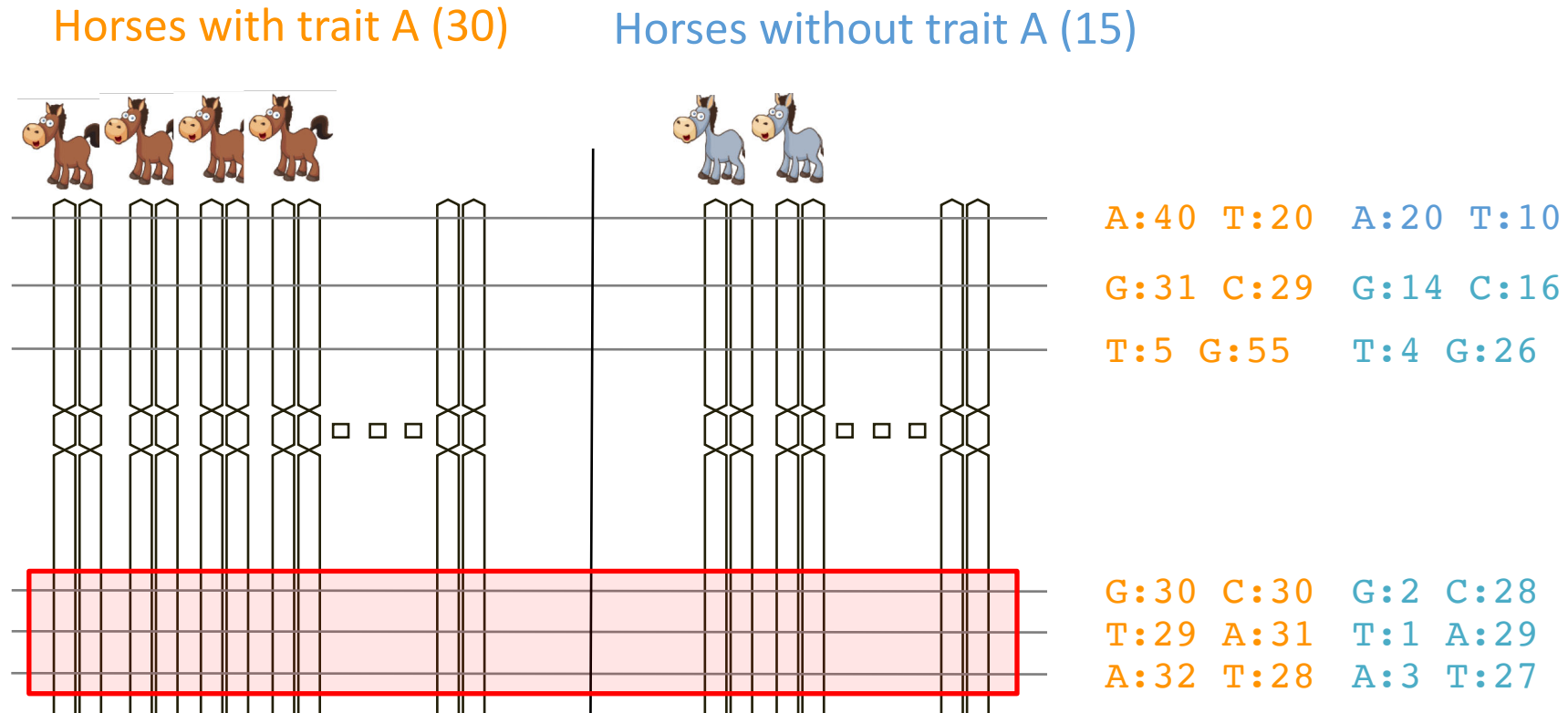
illumina **SPR**

©2013 Neogen Corporation. Neogen and GeneSeek are registered trademarks of Neogen Corporation, Lansing, Mich. All other trademarks are property of their respective companies. A0216-0413

Kit to survey 600,000
SNPs available for horses

- ~400 Thoroughbred horses from JRA
(Currently extended to ~1000 samples)
- Collect blood samples, extract DNA, identify SNPs
- Can we identify genetic variation associated with variation in traits?
- Can we identify genetic variation associated with variation in racing ability?

Genome-wide Association Study (GWAS)



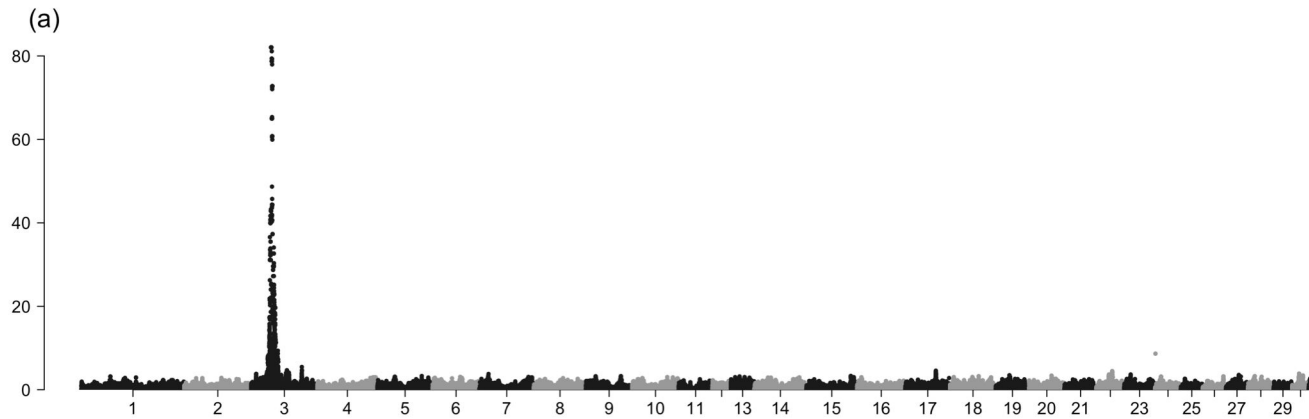
Genetic factor responsible for difference in trait A should be present

Strong association between nucleotide frequency and difference in trait

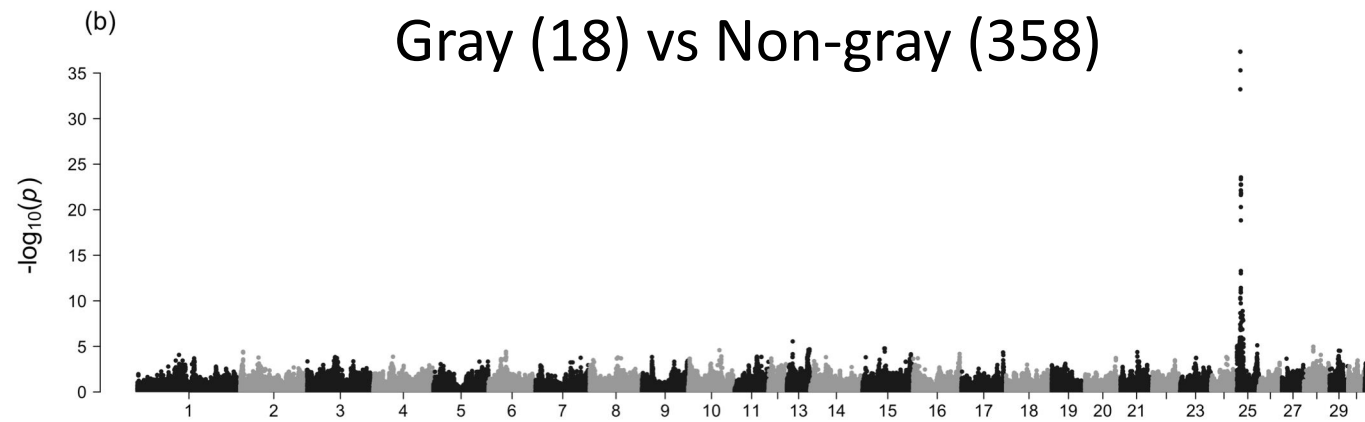
SNPs in neighboring region are "linked"

GWAS works with coat color genes (proof of concept)

Chestnut (94) vs Non-chestnut (258)

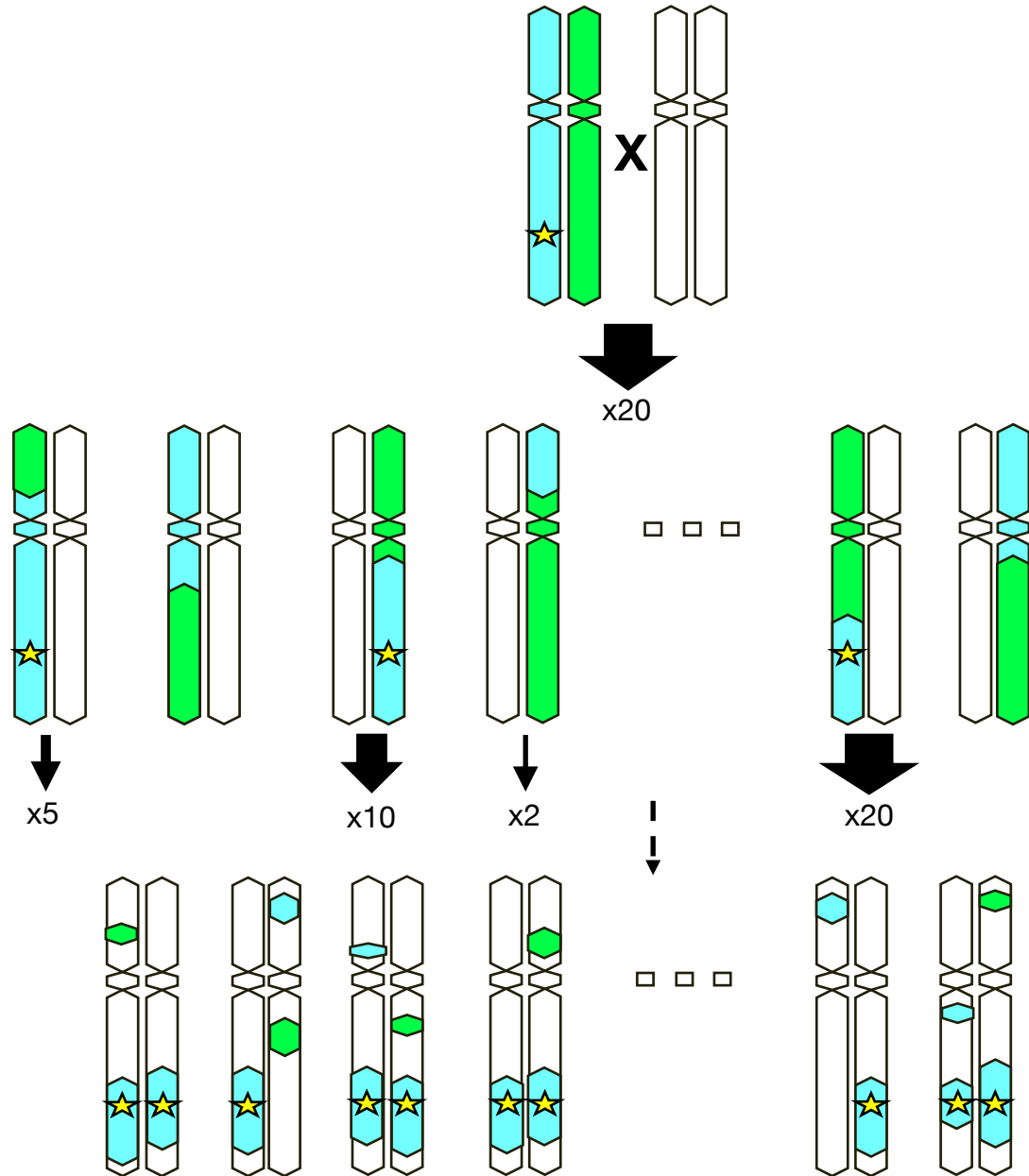


Gray (18) vs Non-gray (358)

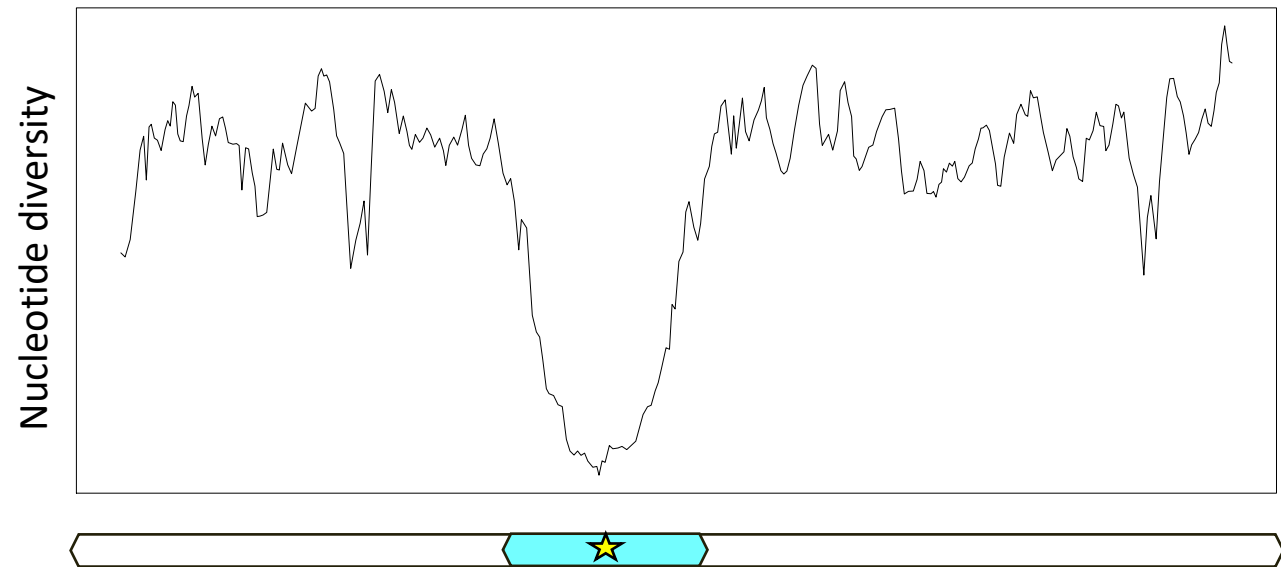


Chromosome

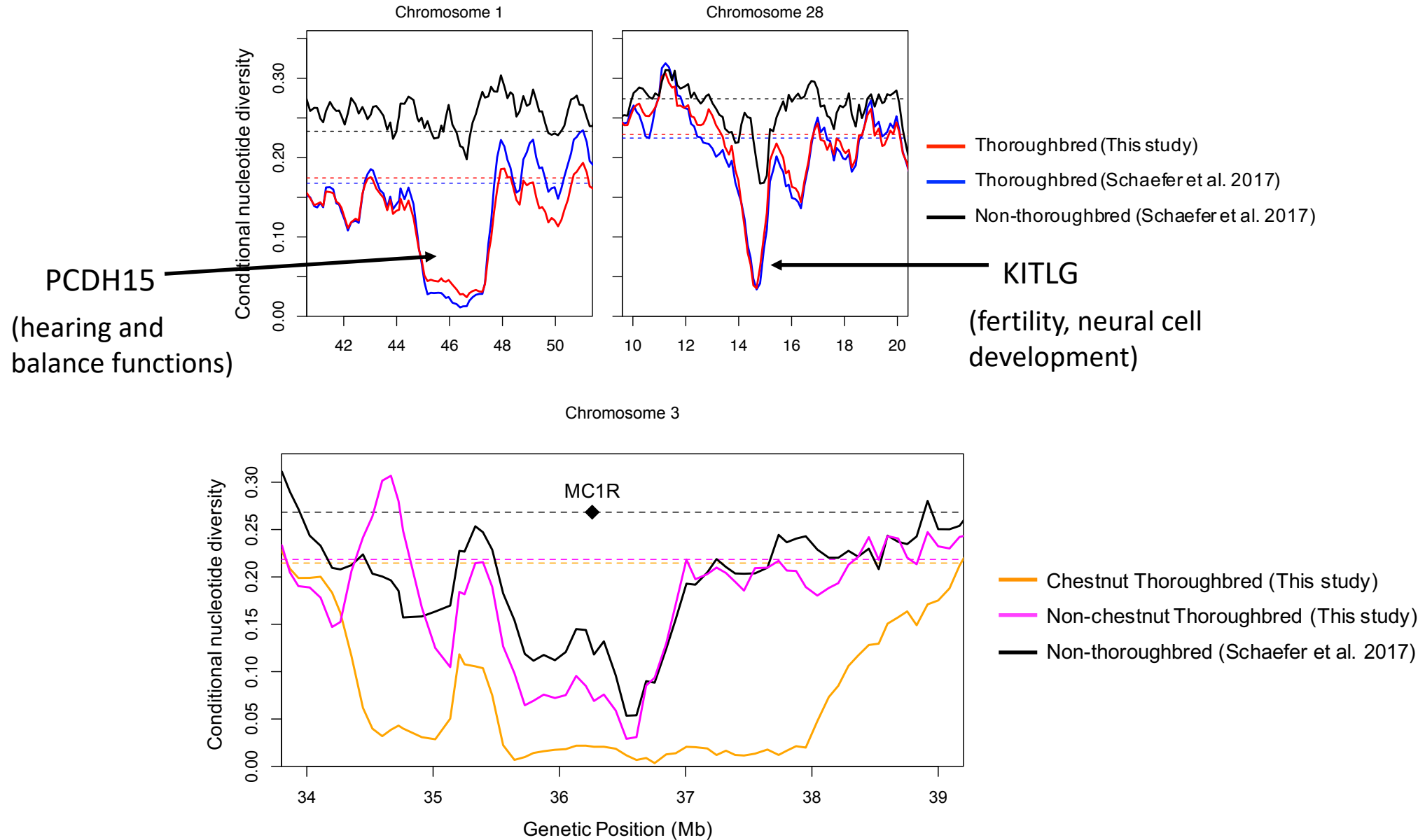
Genomic signatures of selection



- Region containing the mutation has reduced diversity



Regions of reduced diversity regions in Thoroughbreds



Computational approaches to identify important regions

❑ Genome-wide association studies

very powerful when there is a targeted, simple trait

❑ Genomic signature of artificial/natural selection

can use also for complex traits and where there is no specific trait

- Difficult to pinpoint actual causative mutation/gene
- Can identify candidate regions/genes that can be tested experimentally

Summary

- We can explain the process creating the diversity of life in a common mathematical/computational framework
- We can apply the same approaches to study many different species (e.g. horses, buckwheat)
- Computational/mathematical approaches are becoming more important with the increase of data
- More advanced mathematics are needed to go from “sequences” to complex networks, shapes, 3D structures etc?

Acknowledgements

☐ Thoroughbred project

Hideki Innan, Takahiro Sakamoto, Watal M Iwasaki (SOKENDAI), Fumio Sato (JRA Hidaka), Teruaki Tozaki (Lab of Racing Chemistry), Takao Suda (Journalist), etc

☐ Buckwheat project

Yasuo Yasui (Kyoto U), Chengyun Li (Yunnan Agri U), Takanori Ohsako (Kyoto Pref U), Hideki Hirakawa (Kazusa DNA Res), etc