



Seminar „*Genetischer Code und Genomgrammatik*“

Thema: „Statistische Betrachtung des menschlichen Genoms“

Dozent: Prof. Dr. rer. nat. Andreas de Vries
Studierender: Cem Kiyak

25.08.2011

Inhalt

Einleitender Teil:

- Entstehung der „heutigen Zellen“
- Die Zellen mit deren Erbinformationen
- Begrifflichkeiten

Analytischer Teil:

- Statistische Untersuchung

Feststellung

Entstehung der „heutigen Zellen“

- Leben auf der Erde etwa 3,5 Milliarden Jahre alt
- drei große Domänen des Lebens:
Eukaryonten,
Eubakterien und Archaeobakterium
- Annahme:
Verschmelzung des Eubakteriums sowie
Archaeobakteriums für die Entstehung des
Lebens

Die Zellen mit deren Erbinformationen

- Der Körper besteht aus rund **100.000.000.000.000 Zellen**
- Jede biologische Art, jede Spezies, ist in Ihrer Erbinformation einzigartig, sondern sogar jedes Individuum
- Ausnahme: Klon

Arbeiten aus der Praxis: Auf der individuellen Einzigartigkeit basieren forensische Analysen

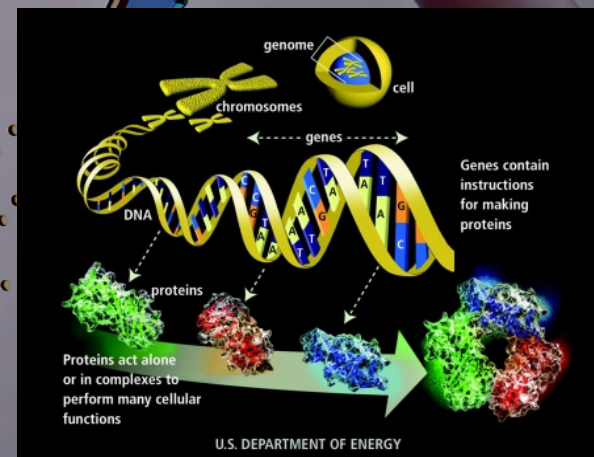
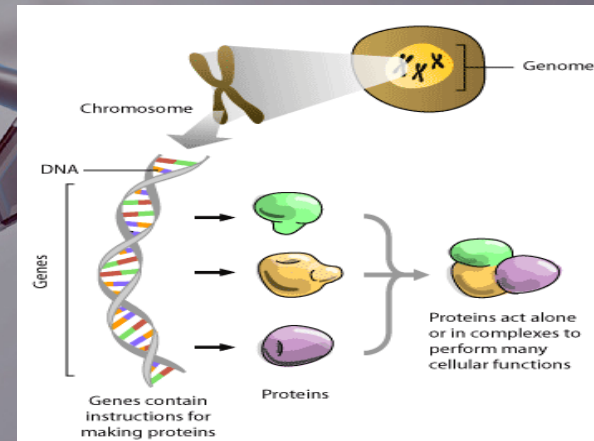
- Die Biologie ist eine Wissenschaft der Grauzonen und Ausnahmen
- Beispiel: **Rote Blutkörperchen**, funktionierende Zellen; aber keine Erbinformation mehr
- Können sich nicht teilen und müssen aus den Stammzellen im Knochenmark gebildet werden
- **Viren** haben Erbinformation, aber sie leben nicht, denn sie brauchen immer lebende Zellen, um Ihre Erbinformation und ihre Hülle zu vermehren.

Die Zellen des Immunsystems

- Erbinformation nicht völlig identisch mit anderen Körperzellen,
- Ihre Erbinformation ist etwas umgestaltet (rekombiniert)

Begrifflichkeiten

- Genetik
- Genom und seine Bedeutung
- Gene der Population
- Gene
- DNA



Charakterisierung eines Gens








- Gen eindeutig charakterisiert
- Durch chemische Bausteine : vier Nukleotiden Adenosin, Cytidin, Guanosin und Thymin
- Baukastenprinzip
- Abfolge dieser Nukleotide, sogenannte Nukleotidsequenz ist ein Gen eindeutig beschrieben
- Drei Milliarden Basenpaare in der DNA des Menschen machen unser Genom aus

Vorgehensweise

- NCBI: National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information aus den USA
- Datenbank der Genomsequenzen per FTP
<http://www.ncbi.nlm.nih.gov/>

Index von ftp://ftp.ncbi.nlm.nih.gov/genomes/H_sapiens/CHR_20/

 [In den übergeordneten Ordner wechseln](#)

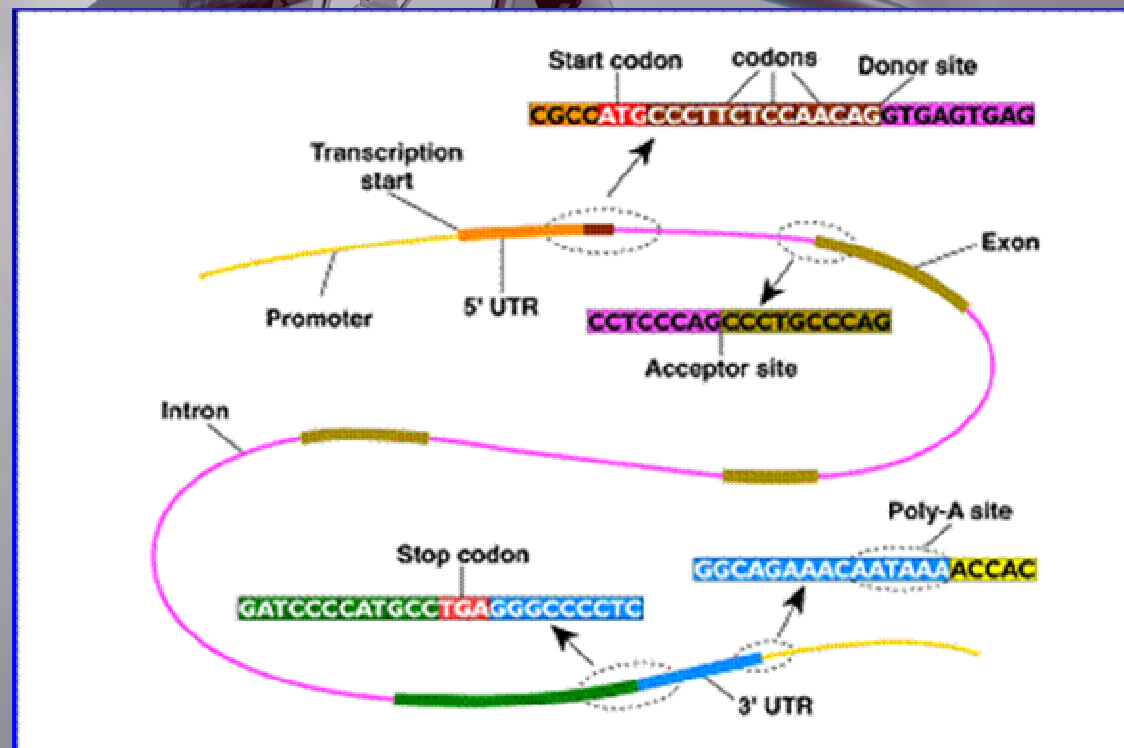
Name	Größe	Zuletzt verändert
 hs_alt_Hs_Celera_chr20.asn.gz	437 KB	22.11.2010 00:00:00
 hs_alt_Hs_Celera_chr20.fa.gz	17408 KB	22.11.2010 00:00:00
 hs_alt_Hs_Celera_chr20.gbk.gz	24236 KB	22.11.2010 00:00:00
 hs_alt_Hs_Celera_chr20.gbs.gz	182 KB	22.11.2010 00:00:00
 hs_alt_Hs_Celera_chr20.mfa.gz	18634 KB	22.11.2010 00:00:00
 hs_alt_HuRef_chr20.asn.gz	541 KB	22.11.2010 00:00:00
 hs_alt_HuRef_chr20.fa.gz	17447 KB	22.11.2010 00:00:00
 hs_alt_HuRef_chr20.gbk.gz	24358 KB	22.11.2010 00:00:00
 hs_alt_HuRef_chr20.gbs.gz	218 KB	22.11.2010 00:00:00
 hs_alt_HuRef_chr20.mfa.gz	18675 KB	22.11.2010 00:00:00
 hs_ref_GRCh37.p2_chr20.asn.gz	376 KB	22.11.2010 00:00:00
 hs_ref_GRCh37.p2_chr20.fa.gz	17533 KB	22.11.2010 00:00:00
 hs_ref_GRCh37.p2_chr20.gbk.gz	24409 KB	22.11.2010 00:00:00
 hs_ref_GRCh37.p2_chr20.gbs.gz	191 KB	22.11.2010 00:00:00
 hs_ref_GRCh37.p2_chr20.mfa.gz	18767 KB	22.11.2010 00:00:00

Java-Tool

- Durch das Java-Tool können nachvollziehbare statistische Auswertungen der FASTA-Dateien bearbeitet werden
- Download unter:
<http://www.math-it.org/java>

25.08.2011

Aufbau des Stranges



Beispiel 1/2 zur Betrachtung des Chromosoms Nr. 20

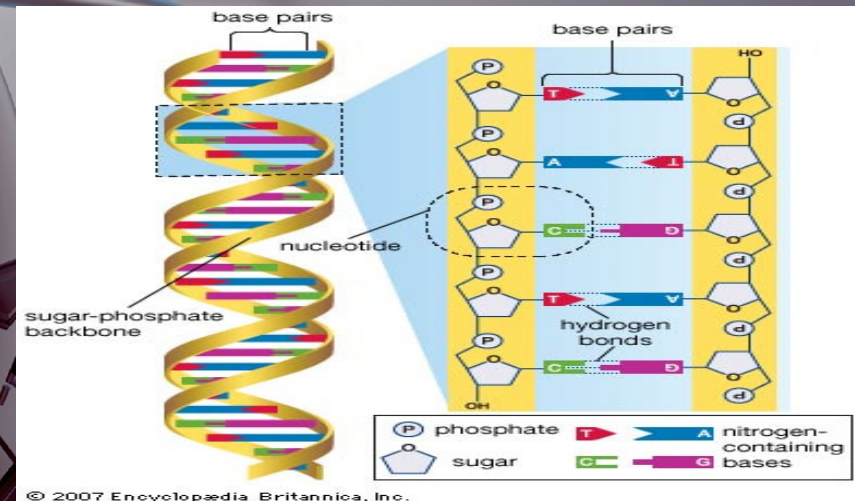
```
CAAATACAAGGGAGAAATTATGATATTCTCAGATAAACAAAAGCTAACAGAGTTCATGACCATTAAACCT
TGCCTACAAGAAATGCTAAAGGGAGTGCTTCAGGTGGAAATGAAAGGATGCTATATGGCAACTCAAAGTT
GTATGAAGAAATAAAGATATCCAGTTATGTGTACATTTTTTAAAAATCATGCTTCAATATCCCACAAACA
AGCTGTGAGCACCCCTGCCCAGGTGACCAGGTAACCAGTATTTTTAAAGCTGGTCCCTGCCACAATTCTCT
TCTTGCTTTCCCTACTGTGACCTAGTGACATTGAAACACTCTAGTAAAAAGTGTTTTTTGCTTGTGTCCTC
CACCTTGACTCCCAGTAAAGGCACCTTGCCCAGGGATCCCATCTCTCTTGCCTACCCATCTGCTTGATTGA
GCCTGCTCCTTGGGAGCTCCTTCCACATGGCTTCTGTCATGGCCTGCCCTACTTCTGTCTCTAGGAACTG
TGAGTATGCCTCTTCATATGAACTTTTCATGGCCATGTTAGAGTGATCCTCTAAGATTCAACCAGCAGGA
GTCACTCTAACATTTTCTATAGTAATGGCAATGAGGATGGGATCATTCTAAACTCAGGTGGAAAGTTCTT
AGGGAAATGCCTTTGACTCATGACTTACTGATTTGCTACTCTGATTGGCTCCACAATTTGAGAAGGCTTTT
TACCCTGATGACTTTCTTGTGCTATGCTATGTACAGCTTTGTGTTGCCTTTTGGAGTGCTGCCATGGCTC
TCCATCCTGAAGCGAGGATCCTGTTCCAATATCAATAATGACATCCACCTTTCTCCAAGAGCACAAAGATG
ATATAACATTGGTAACACCTCTTCAAAGAAGATTCATGTGATATGCATATTAATGTCACATAATCCAGCA
CCTGTAAGTTTTTTTTTTGGGATGCCTGTGGCAACATATTCTTCATTTACAAAATTTTACTTACAAAATTA
TCAAACAGACACTCCTGTTAGTGTCATTAATGAATTAACAATAAATATCTCTCCTCTGTAAAGACTTTGG
CTGTGCTTTTCATGCCTCTTGGAGTCTCTAGACCACCTCATGATGGCTATATATTTTCTCAGGGCCTTGAT
GCAAAAAATAATGCTCTTTATGGCCTACGCTGTACTCCATTAAAAACAACCTGACTGGTTACCCACTAGAATC
TCCTGATTTGGATTCCAGATTTGAAGATAATACCCTCCTGTTTCAATCATTAAACCATAAAATTGCCCCACA
GGCCAGGCTCTGCTCATCTGAAATGACACCCATCGGAATTTAAACTTCCACTGCAACACAAAATTAATAAT
TGCTGCCAGGTACAGCTGCACACCTTTCTGTTGACAGTTGGCATTATTACTACTTGGTTTGTGATGCAA
CAATCCTCAACCATAACTAGCAAACTGAATTGAACAGCATATTAGAAGAATTATACAACACGAGCAAGT
GATATTTATTTCTGGAATACAAGGATGTTTCCATGCACAAATGACAAACAATGTAATATATCATATTTAT
AGAAGAAAGACAAGAAACACATGATTAGCTCATTTGATACAGAAGACGATTTGACAAAATTTGGGGATCTT
TTTATGACAAAAACCATCAATAAAAAGGAATATATGAAAACCTTCTATAAAAATAAATAAAGTCGTAAACACA
CACTTGATTTTCACTTATCGCTGTCTAATAATACTATTTTATTTGCCTATTTATCATGTGTATTGATTACT
GTCATTCTCCCCAGTGGAGATAAGCAACAAAAGGGCAACTACTTTTTGTCTTTTATTCAACAATGTTCTT
TGTTAGGTCCTGGCACATAGCAGGTATTCATCAAAATATGTGTATAATGTGTGAATTGCTTAGTGGAACAC
```

Beispiel 2/2 zur Betrachtung des Chromosoms Nr. 20

```
AGGATCTCTGCATGGGGGCAGGAGTCTCTGCGGGGAGGACTTTTtagggggtaaacaggatctgtgcagcgt
GCAGGATCTTTGcAGTGGGGCCAGGATCTGTGCAGTGGGTAGGATCTACATTGGGAGTCAGTGTCTGTTC
AAGGGAGTCGGGGGACAGGACCCGTCCAGGGAGGGCAGGATCTGTACAGAGGGACCAAGATCTACGCACG
GGGATAGCATCTATGTGGGGTGGGGGCACGGTCTGTGCCAGGTTTTGGGGGC ACTAGATCTTTGcAGGTG
TGGCCAGGATCTGTGC AAGGGGTAGGATCTGTGCAGAGTGGCCAGGATCTGTGCAGAGTGGCCAGGATCT
GTGCAGGGAGGGCAGGATCTGTGCAGAGTGGCCAGGATCTGTGCAGAGTGGCCAGGATCTGTGCAGAGTG
GCCAGGATCTGTGCAGGGAGGCCAGGATCTGTGCAGAGTGGCC ANNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
```

Analyse des menschlichen Genoms nach der Länge der Nukleotidbasen

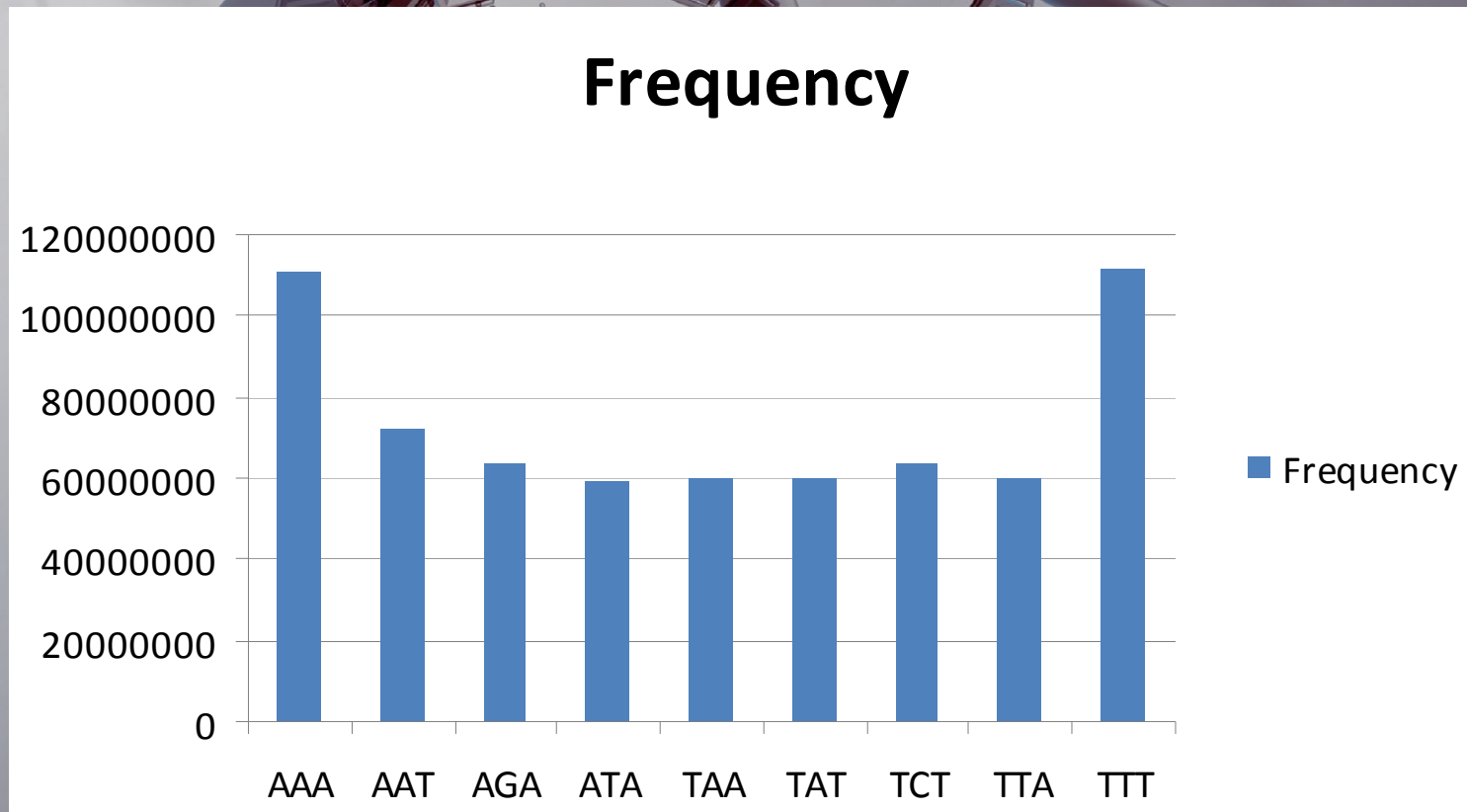
- Nukleotidbasen der Länge 3 (AAA,AAT,AGA,...)
- Nukleotidbasen der Länge 6 (AAAAAA,...)
- Nukleotidbasen der Länge 9 (AAAAAAAAAA,...)



1	GG	AGG	ATTT	GG	AA	ATT	TAT	TTT	AA	TT	24
2	5	CC	ATT	AA	TAT	TAT	AGG	ATC	ACC	TGA	48
4	9	A	TAG	CATT	CCC	AC	GAT	TAA	AA	TAA	72
7	3	A	TTA	GAT	TTT	TG	A	CT	TTT	TAC	96
9	7	T	CAT	TAT	TTT	A	TAT	TAA	AT	TTA	120
1	2	1	AA	TAT	AT	TTT	A	TAC	CTA	AT	144
1	4	5	A	C	AGG	AT	G	A	A	CTT	168
1	6	9	T	T	A	T	C	T	T	C	192
1	9	3	T	C	A	C	T	T	C	A	216
2	1	7	T	T	T	C	A	T	T	T	240
2	4	1	T	C	C	T	C	T	A	T	264
2	6	5	T	T	A	T	T	G	T	A	288
2	8	9	A	A	A	A	T	T	T	T	312
3	1	3	C	A	A	A	T	T	A	T	336
3	3	7	G	T	A	T	G	T	A	T	360

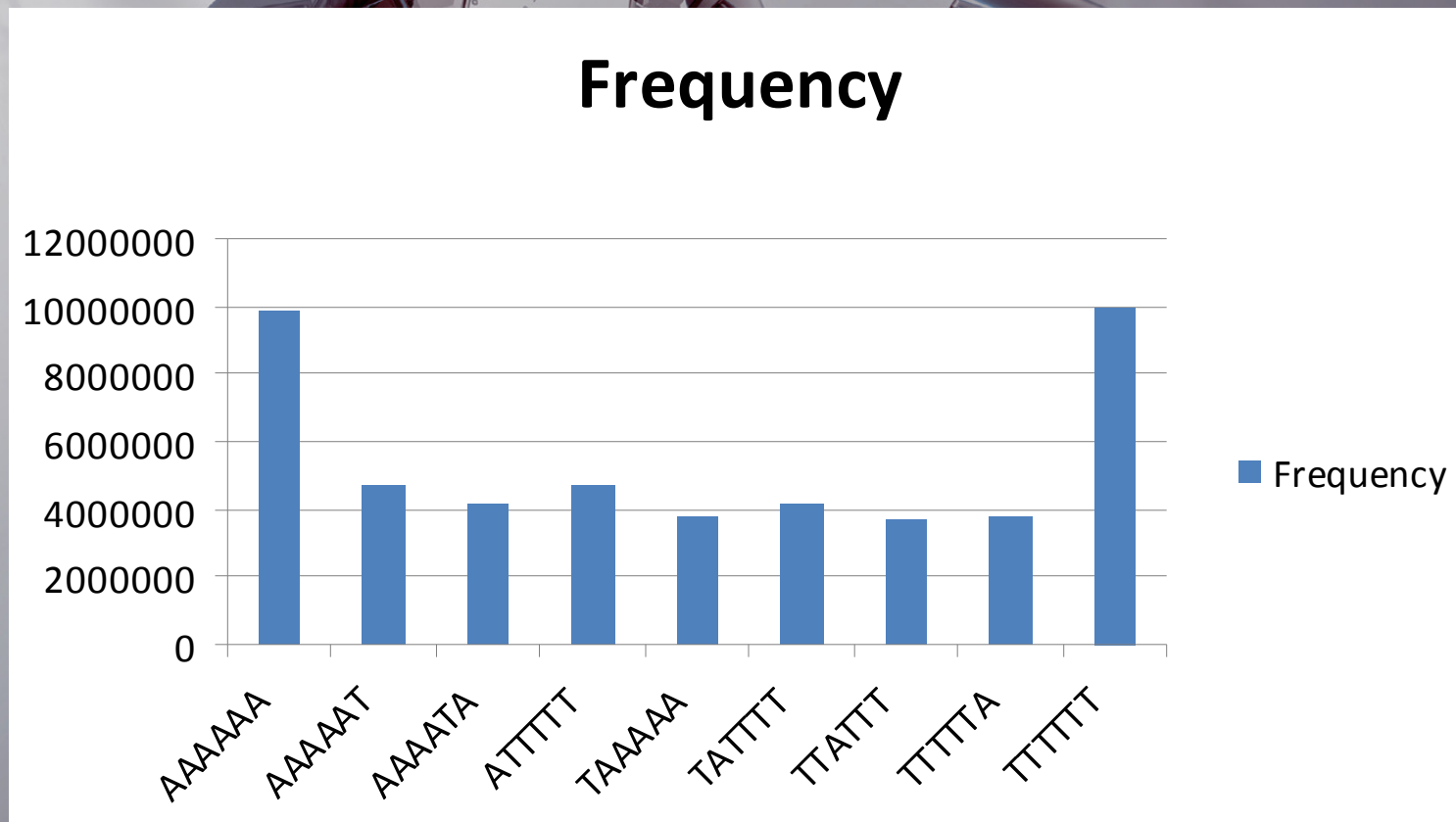
Darstellung der Häufigkeiten

- Nukleotidbasen der Länge 3



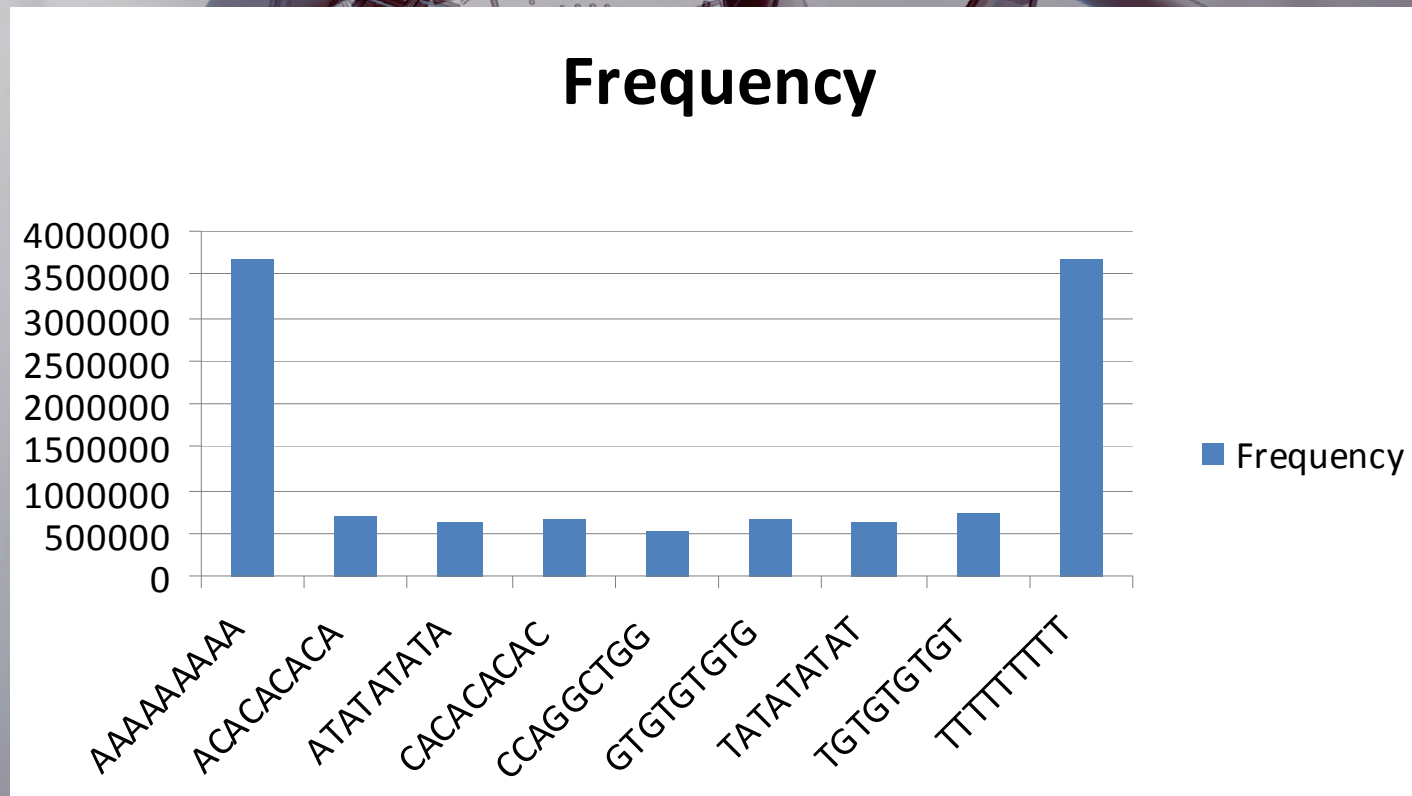
Darstellung der Häufigkeiten

- Darstellung der Häufigkeiten der Länge 6



Darstellung der Häufigkeiten

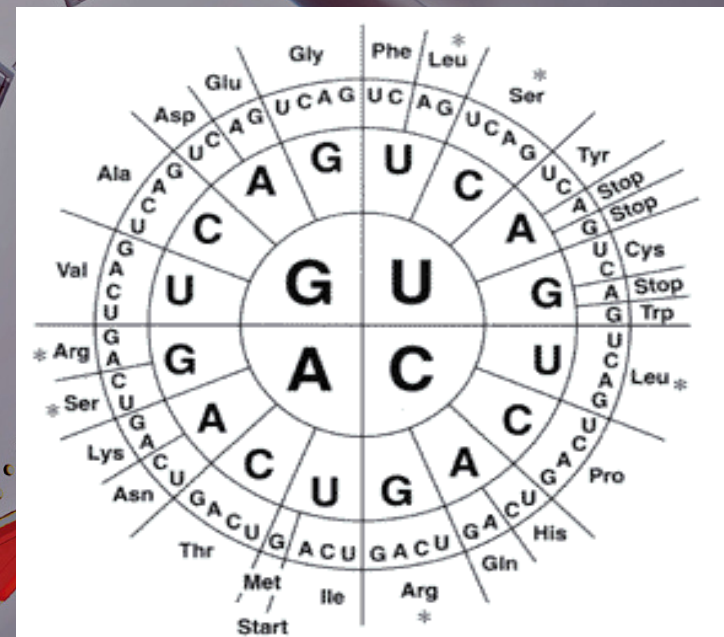
- Darstellung der Häufigkeiten der Länge 9



Die nicht genutzten Kombinationen der Länge 9 des menschlichen Genoms

Die 6 häufigsten davon: Codierung der Nucleotide:

- CGTCGATCG 9
- CGACGTACG 7
- CGATACGCG 7
- CGTAACGCG 7
- TACGCGCGA 7
- TCGCGTCGA 7



Aufspaltung in Condons sowie Aufschlüsselung (Erkenntnis)

Aufspaltung des Genoms in seine Basen (3 davon jeweils = Codons bzw. Triplets) und auffinden seiner Aminosäuren:

- 1) CGT = Arg
 - 2) CGA = Arg
 - 3) CGA = Arg
 - 4) CGT = Arg
 - 5) TAC = Tyr
 - 6) TCG = Ser
- CGA = Arg
CGT = Arg
TAC = Tyr
AAC = Asn
GCG = Ala
CGT = Arg
- TCG = Ser
ACG = Thr
GCG = Ala
GCG = Ala
CGA = Arg
CGA = Arg

Nicht vorhandene Kombinationen der Länge 12 (Erkenntnis)

	A	B	C	D	E	F	G	H	I
1	12-letter Words omitted in all human chromosomes								
28									
29	AAAAACGCGCGA	AAATCGGCGACG	AACCGGATGACG	AACGAGTGTTCG	AACGCGGGCGCAT	AACGGTCAGTCG	AACTACGTCCGC	AAGCGCTAGCGC	AATAAGTTCGCG
30	AAAAACGCGCTA	AAATCGGCGCGC	AACCGGCACGTA	AACGAGTTCGAT	AACGCGGGCGCTA	AACGGTCATGCG	AACTACGTCCGC	AAGCGCTATCCG	AATAATCCGCCG
31	AAAAATCGCGAC	AAATCGGCGCTT	AACCGGCCCGCT	AACGAGTTTACG	AACGCGGGCGTT	AACGGTCATTCG	AACTACGTCCGG	AAGCGCTATGCG	AATAATCCGCCG
32	AAAAATCGCGCG	AAATCGGCGTAA	AACCGGCCCGTA	AACGATAACCGC	AACGCGGGCGGAT	AACGGTCCACGA	AACTACTCGTCG	AAGCGCTCGACG	AATAATCGACGG
33	AAAACCGACGCG	AAATCGGCGTCG	AACCGGCCGAAT	AACGATAATCCG	AACGCGGGCTAA	AACGGTCCCGAC	AACTACTGCGCG	AAGCGCTTACGC	AATAATCGCCG
34	AAAACCGCGCTA	AAATCGGCTCGT	AACCGGCCGCGT	AACGATACACGG	AACGCGGGCTAC	AACGGTCCGACC	AACTAGACGCGG	AAGCGCTTATCG	AATAATCGTCGG
35	AAAACCGGTTTCG	AAATCGGGCGAT	AACCGGCCGAACG	AACGATACATCG	AACGCGGGCTACG	AACGGTCCGACG	AACTAGACGCGT	AAGCGGACAACG	AATACACCGCGA
36	AAAACCGTCGCG	AAATCGGGTACG	AACCGGCCGAATT	AACGATACCGCC	AACGCGGGCTCGA	AACGGTCCGCGG	AACTAGCCGCGG	AAGCGGACCGGT	AATACACCGGCG
37	AAAACCGTCGCT	AAATCGGGTCGA	AACCGGCCGACAT	AACGATACCGCG	AACGCGGGTTCG	AACGGTCCGGTA	AACTAGCCGCGA	AAGCGGACGATA	AATACACCGTCG
38	AAAACGACCGGT	AAATCGGTACGA	AACCGGCCGATTA	AACGATACCGGT	AACGCGGGTAGT	AACGGTCCGGTT	AACTAGCCGTCG	AAGCGGACGTAC	AATACACGGCGT
39	AAAACGACGCGC	AAATCGGTGCGC	AACCGGCCGCAAT	AACGATACGACC	AACGCGGGTATA	AACGGTCCGTAA	AACTAGCGCCGT	AAGCGGATACGT	AATACCACGTG
40	AAAACGATCGGG	AAATCGGTTGCG	AACCGGCCGCAAT	AACGATACGACT	AACGCGGGTCTA	AACGGTCCGTCA	AACTAGCGCGAA	AAGCGGATTCG	AATACCAGTCCG
41	AAAACGCCGATC	AAATCGTAACCG	AACCGGCCGCTA	AACGATACGAGC	AACGCGGGTTAC	AACGGTCCGTCC	AACTAGCGCGCT	AAGCGGATTCGG	AATACCAGCCGA
42	AAAACGCCGTAC	AAATCGTACCCG	AACCGGCCGCGCA	AACGATACGCGG	AACGCGGGTACGA	AACGGTCCGTCCG	AACTAGCGCGTT	AAGCGGCCCGTA	AATACCAGCGAT
43	AAAACGCGACCG	AAATCGTACGGC	AACCGGCCGCGGT	AACGATACGCGA	AACGCGGGTACGG	AACGGTCCCTATC	AACTAGCGCTCG	AAGCGGGCATCT	AATACCAGCGCC
44	AAAACGCGGACG	AAATCGTCGACG	AACCGGCCGCTAA	AACGATACGCGC	AACGCGGGTACGT	AACGGTCCGAACC	AACTAGCGTACG	AAGCGGCCGCTA	AATACCAGTCGA
45	AAAACGCGATCG	AAATCGTCGCAA	AACCGGCCGCGCT	AACGATACGCGG	AACGCGGGTACTC	AACGGTCCGAATC	AACTAGCGTCA	AAGCGGCCGACG	AATACCAGTCGG
46	AAAACGCGCGAT	AAATCGTCGCAC	AACCGGCCGTTA	AACGATACGCGG	AACGCGGGTAGAA	AACGGTCCGACAT	AACTAGCGTCCG	AAGCGGCCGATA	AATACCCTCGCG
47	AAAACGCGCGTA	AAATCGTCGCGC	AACCGGCCGTACG	AACGATACGGCT	AACGCGGGTATCA	AACGGTCCGACCG	AACTAGTACGGG	AAGCGGCTAACG	AATACCAGCCGA
48	AAAACGCGTCGG	AAATCGTCGCGT	AACCGGCCGTAGT	AACGATACGTTT	AACGCGGGTATT	AACGGTCCGACGA	AACTAGTCCGGA	AAGCGGCTATCG	AATACCAGCCGG
49	AAAACGCGTTTCG	AAATCGTCGCTG	AACCGGCCGTCAA	AACGATAGCCGT	AACGCGGGTCAAC	AACGGTCCGACGC	AACTAGTCCGCG	AAGCGGGTATCG	AATACCAGCGAC
50	AAAACGCTTCGC	AAATCGTCGGGG	AACCGGCCGTCGT	AACGATAGCGAC	AACGCGGGTCAAT	AACGGTCCGACTC	AACTAGTCCGCT	AAGCGGGTCCGG	AATACCAGCGCG
51	AAAACGGCGCCG	AAATCGTCGGTT	AACCGGCCGTGTA	AACGATAGCGCC	AACGCGGGTCCAA	AACGGTCCGACTG	AACTAGTCCGCG	AAGCGGGTCCGAA	AATACCAGTCCG
52	AAAACGGCGCCA	AAATCGTCGGCC	AACCGGCCGTTAA	AACGATAGCGCC	AACGCGGGTCCCA	AACGGTCCGACCC	AACTACTCCGCG	AAGCGGTAATCC	AATACCAGCCCGA

Anzahl: 128 000 Kombinationen nicht vorhanden

Was wollen wir bewirken?

- Defekte Gene „reparieren“
- Änderung der Gene um neues oder anderes Leben zu erzeugen
- Krankheiten besser verstehen und evt. Lösungen dazu zu finden



Feststellung



- Intensive weitere Analysen von großer Bedeutung, denn hier ist der Weg das Ziel!
- Verstehen der Codon-Darstellung die Aufschluss auf eine verständliche Sprache ergeben sollen
- Die Chancen könnten gut stehen durch statistische Berechnungen mit durchaus umfangreicheren Untersuchungen und leistungsstarken Rechnern und Geduld ein Ergebnis zu erhalten
- Sind Grammatikalische Strukturen zu erkennen?



Vielen Dank für Ihre Aufmerksamkeit!

25.08.2011

21