

WS2017/2018 MSc Modul 7A

# **„Genomforschung und Sequenzanalyse**

- Einführung in Methoden der Bioinformatik- “

AG Hankeln

---



## **Trio-Analyse (Exome-Seq)**

**Mutationen mit  
Krankheitsrelevanz  
im Humangenom**

# More than 3800 genes are (known to be) involved in human disease...

**Online Mendelian  
Inheritance of Man  
= catalogue of human  
genes and diseases**

## OMIM Gene Map Statistics

OMIM Morbid Map Scorecard (Updated November 16th, 2017) :

Total number of phenotypes* for which the molecular basis is known	6,112
Total number of genes with phenotype-causing mutation	3,844
* Phenotypes include (1) single-gene mendelian disorders and traits; (2) susceptibilities to cancer and complex disease (e.g., BRCA1 and familial breast-ovarian cancer susceptibility, <a href="#">113705.0001</a> , and CFH and macular degeneration, <a href="#">134370.0008</a> ); (3) variations that lead to abnormal but benign laboratory test values ("nondiseases") and blood groups (e.g., lactate dehydrogenase B deficiency, <a href="#">150100.0001</a> and ABO blood group system, <a href="#">110300.0001</a> ); and (4) select somatic cell genetic disease (e.g., GNAS and McCune-Albright syndrome, <a href="#">139320.0008</a> and IDH1 and glioblastoma multiforme, <a href="#">147700.0001</a> .)	

Distribution of Phenotypes across Genes (Updated November 16th, 2017) :

Number of genes with 1 phenotype	2,631
Number of genes with 2 phenotypes	721
Number of genes with 3 phenotypes	259
Number of genes with 4+ phenotypes	233

Dissected OMIM Morbid Map Scorecard (Updated November 16th, 2017) :

Class of phenotype	Phenotype	Gene *
Single gene disorders and traits	5,072	3,464
Susceptibility to complex disease or infection	698	499
"Nondiseases"	145	115
Somatic cell genetic disease	212	121
*Some genes may be counted more than once because mutations in a gene may cause more than one phenotype and the phenotypes may be of different classes (e.g., activating somatic BRAF mutation underlying cancer, <a href="#">164757.0001</a> . and germline BRAF mutation in Noonan syndrome, <a href="#">164757.0022</a> .)		

+141900  
Table of Contents

Title

Gene-Phenotype  
Relationships

Clinical Synopsis

Text

Description

Gene Structure

Mapping

Gene Function

Biochemical Features

Molecular Genetics

Animal Model

History

Allelic Variants

Table View

See Also

References

Contributors

Creation Date

Edit History

+ 141900

## HEMOGLOBIN--BETA LOCUS; HBB

Other entities represented in this entry:

**METHEMOGLOBINEMIA, BETA-GLOBIN TYPE, INCLUDED**  
**ERYTHREMIA, BETA-GLOBIN TYPE, INCLUDED**

*HGNC Approved Gene Symbol:* **HBB**

*Cytogenetic location:* **11p15.4**    *Genomic coordinates (GRCh38):* **11:5,225,465-5,227,070** (from NCBI)

### Gene-Phenotype Relationships

Location	Phenotype	Phenotype MIM number	Inheritance	Phenotype mapping key
11p15.4	Delta-beta thalassemia	141749	AD	3
	Erythremias, beta-			3
	Heinz body anemias, beta-	140700	AD	3
	Hereditary persistence of fetal hemoglobin	141749	AD	3
	Methemoglobinemias, beta-			3
	Sickle cell anemia	603903	AR	3
	Thalassemia-beta, dominant inclusion-body	603902		3
	Thalassemias, beta-	613985		3
	{Malaria, resistance to}	611162		3

Clinical Synopsis ▾

### TEXT

#### ▼ Description

The alpha (HBA1, 141800; HBA2, 141850) and beta (HBB) loci determine the structure of the 2 types of polypeptide chains in adult hemoglobin, HbA. Mutant beta globin that sickles causes sickle cell anemia (603903). Absence of beta chain causes beta-zero-thalassemia. Reduced amounts of detectable beta globin causes beta-plus-thalassemia. For clinical purposes, beta-thalassemia (613985) is divided into thalassemia major (transfusion dependent), thalassemia intermedia (of intermediate severity), and thalassemia minor (asymptomatic).

#### ▼ External Links

► Genome

► DNA

► Protein

► Gene Info

#### ▼ Clinical Resources

Clinical Trials  
Gene Tests  
► Genetics Home  
Reference  
GTR  
Newborn Screening  
GARD

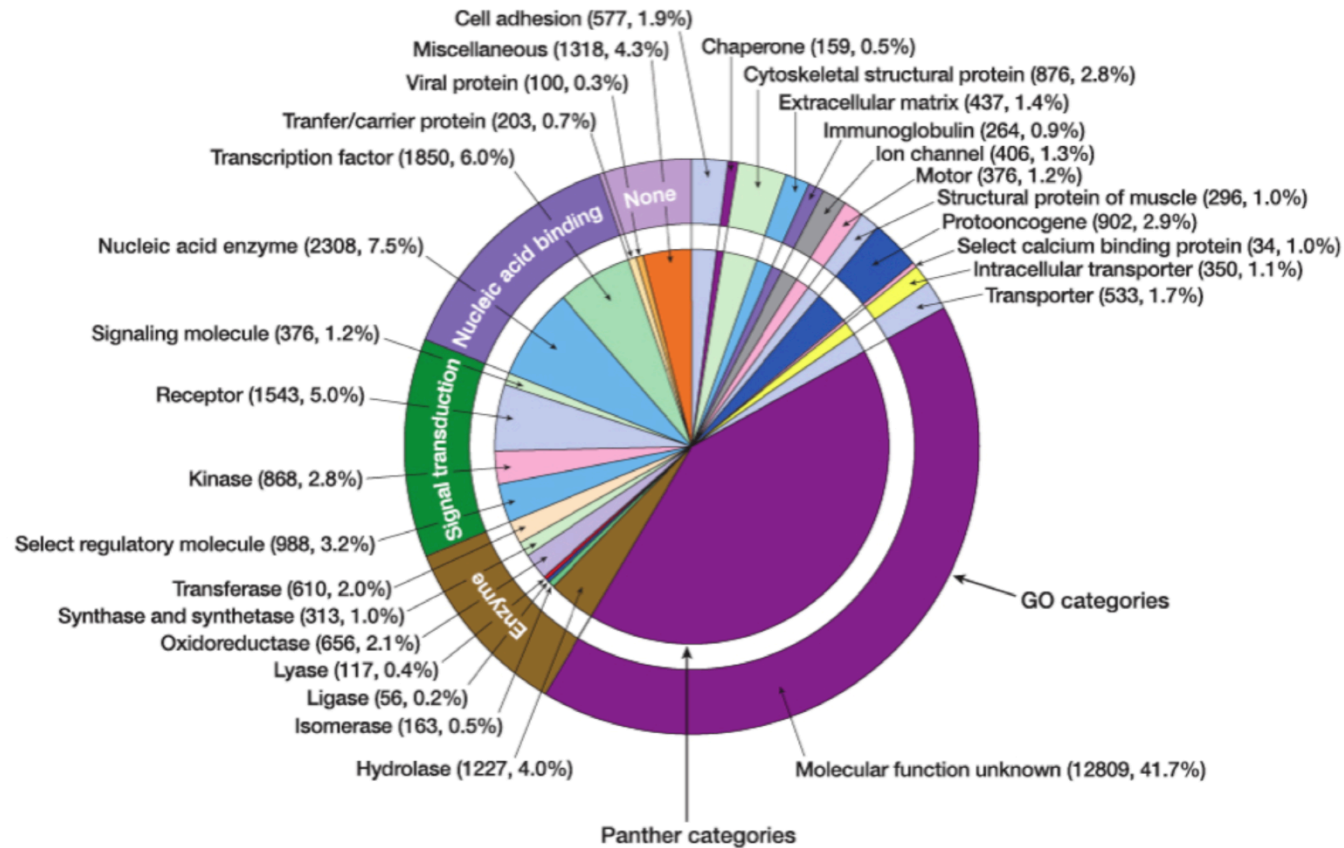
► Variation

► Animal Models

► Cell Lines

► Cellular Pathways

# Functions of human genes



**Figure 12.20: A preliminary functional classification of human polypeptide-encoding genes.**

Known or predicted functions for 26 383 human polypeptide-encoding genes. Classification is according to the GO molecular function categories as shown in the outer circle (Gene Ontology classification – see Section 8.3.6) or to Celera's Panther molecular function categories (inner circle). Reproduced from Venter *et al.* (2001) *Science* **291**, 1304–1351, with permission from the American Association for the Advancement of Science.



# Nomenklatur bei Genveränderungen

## SNV

Single Nucleotide *Variant*

Persönliche Variation des  
Genoms

## SNP

Single Nucleotide *Polymorphism*

SNV, der in mindestens 1% der  
untersuchten Population  
gefunden werden konnte

## CNV

Copy number variation

# SNVs

## Single Nucleotide Variants

..AC <b>G</b> GC..	..AC <b>T</b> GC..
..TG <b>C</b> CG..	..TG <b>A</b> CG..

Zwei verglichene Genome unterscheiden sich im Mittel alle **1000 Bp**.

Bei Vergleich aller Genome weltweit schätzt man **9-12 Mio. SNVs (also 1/300 Bp)**.

Aus diesen (und weiteren) **Unterschieden** resultiert unsere **Individualität**

# 1000 Genomes

A Deep Catalog of Human Genetic Variation



[Home](#) [About](#) [Data](#) [Analysis](#) [Participants](#) [Contact](#) [Browser](#) [Wiki](#) [FTP search](#)

Search

## LATEST ANNOUNCEMENTS

WEDNESDAY SEPTEMBER 30, 2015

### A global reference for human genetic variation

The Phase 3 publication, [A global reference for human genetic variation](#) and the Phase 3 Structural variation publication, [An integrated map of structural variation in 2,504 human genomes](#) are now available from [Nature](#) alongside a [celebration of 25 years of the Human Genome Project](#)

The variants from the Phase 3 analysis are available in <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/> and extended information about the SV dataset can be found in [ftp/phase3/integrated\\_sv\\_map/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/ftp/phase3/integrated_sv_map/).

Both these papers are open access and should be free for everyone to read and download.

If you have any questions about the data these papers are based on or how to access it please email [info@1000genomes.org](mailto:info@1000genomes.org)

<http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>

### Recent project announcements

FRIDAY NOVEMBER 27, 2015

#### Changes to the 1000 Genomes Globus endpoint

EMBL-EBI has recently rearranged its Globus hosted endpoints.

## NAVIGATION

- [Frequently Asked Questions](#)

## LINKS



[All Project Announcements](#)



[Sample and Project Information](#)



[Media Archive](#)

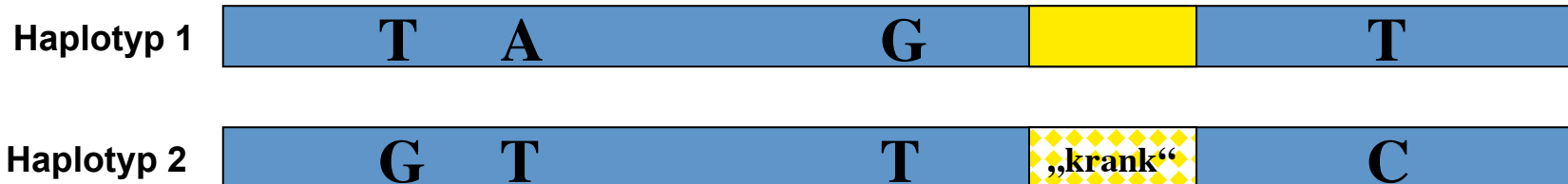


[Find the 1000 Genomes Project Publications](#)

# SNPs, Haplotypen und Erkrankungen

**ACHTUNG:** nur selten ist ein SNP die direkte genetische Ursache für eine Erkrankung oder Prädisposition!

SNPs sind vielmehr als **genetische Marker** zu sehen, die einen „**Haplotyp**“ definieren und mit einem unbekannten Krankheitsallel **gekoppelt** vererbt werden



unbekanntes Gen

# Genetische Krankheiten

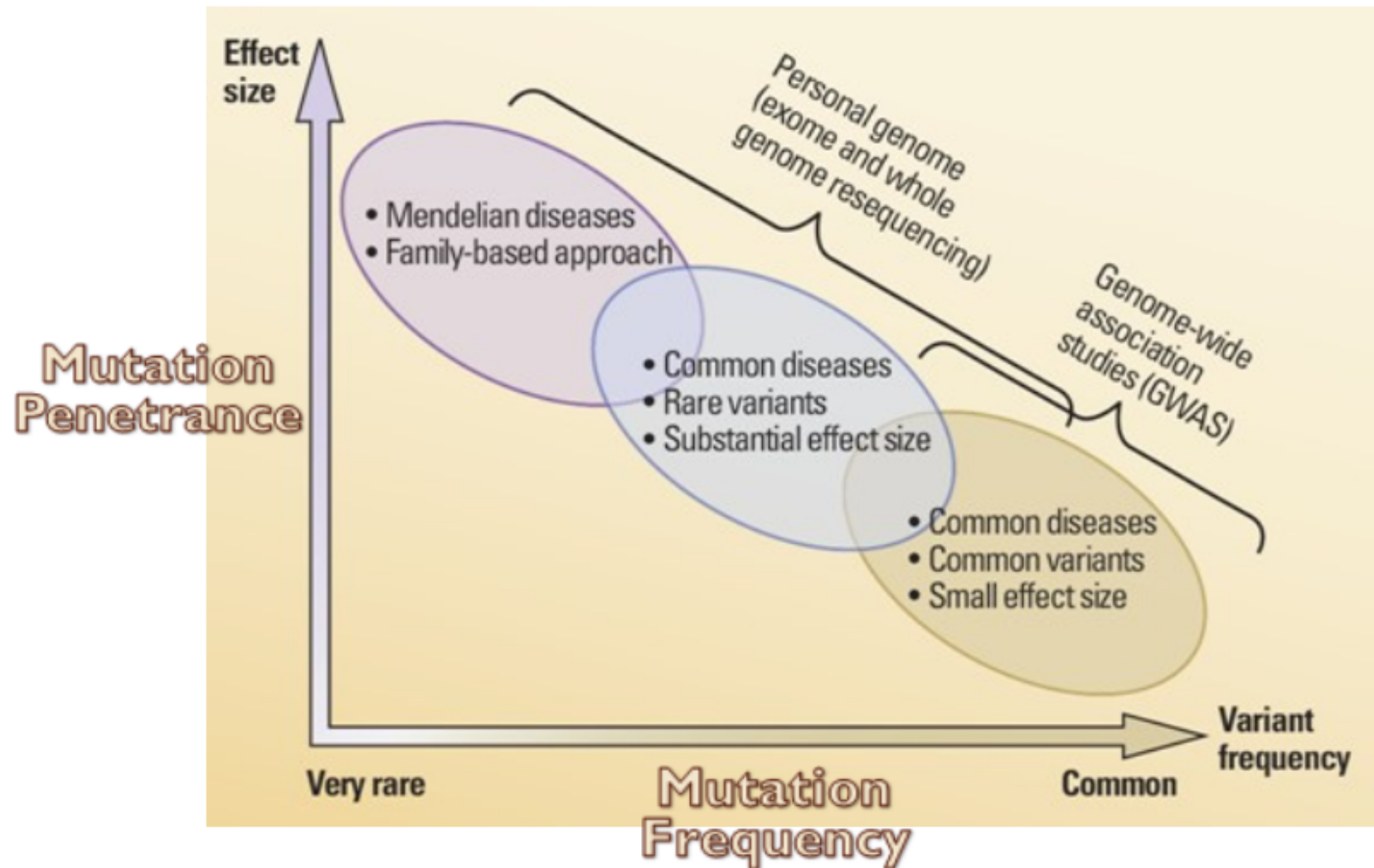
- **Mendel'sche Krankheiten**

- (meist) monogen
- Hohe Penetranz
- Beispiele: Sichelzellenanämie (HBB), zystische Fibrose (CFTR), Huntington

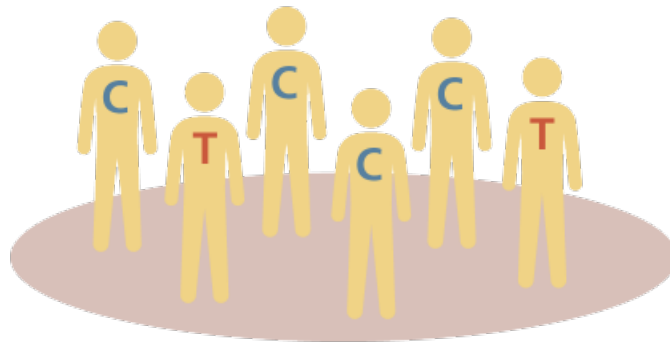
- **Komplexe Krankheiten**

- Polygen, d.h. durch mehrere Gene bedingt
- Niedrige Penetranz
- Multifaktoriell: können durch Umwelt/Ernährung beeinflusst werden
- Beispiele: Herz-Kreislaufkrankungen, Diabetes, Übergewicht, Demenz, Krebs, Skelett

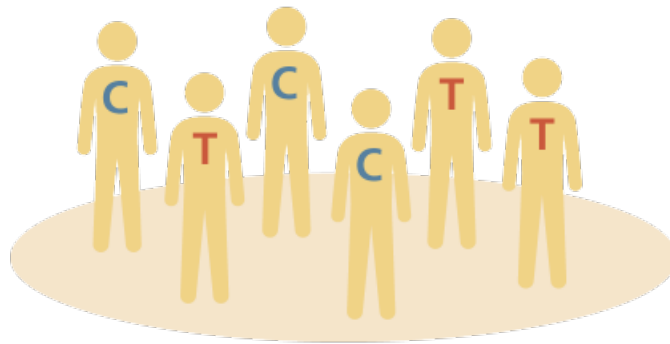
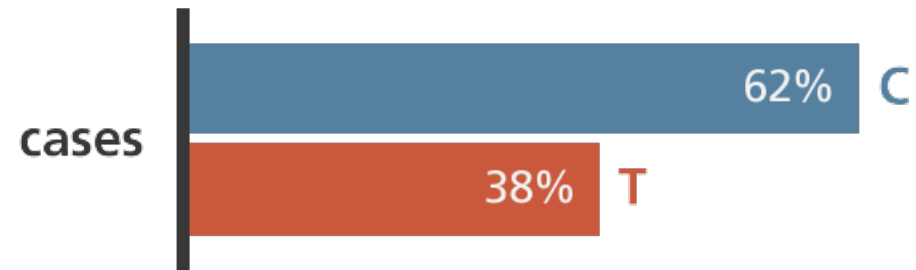
# SNPs und Erkrankungen



# Genome-wide association studies (GWAS)



**cases (n=1,000)**  
people with heart disease

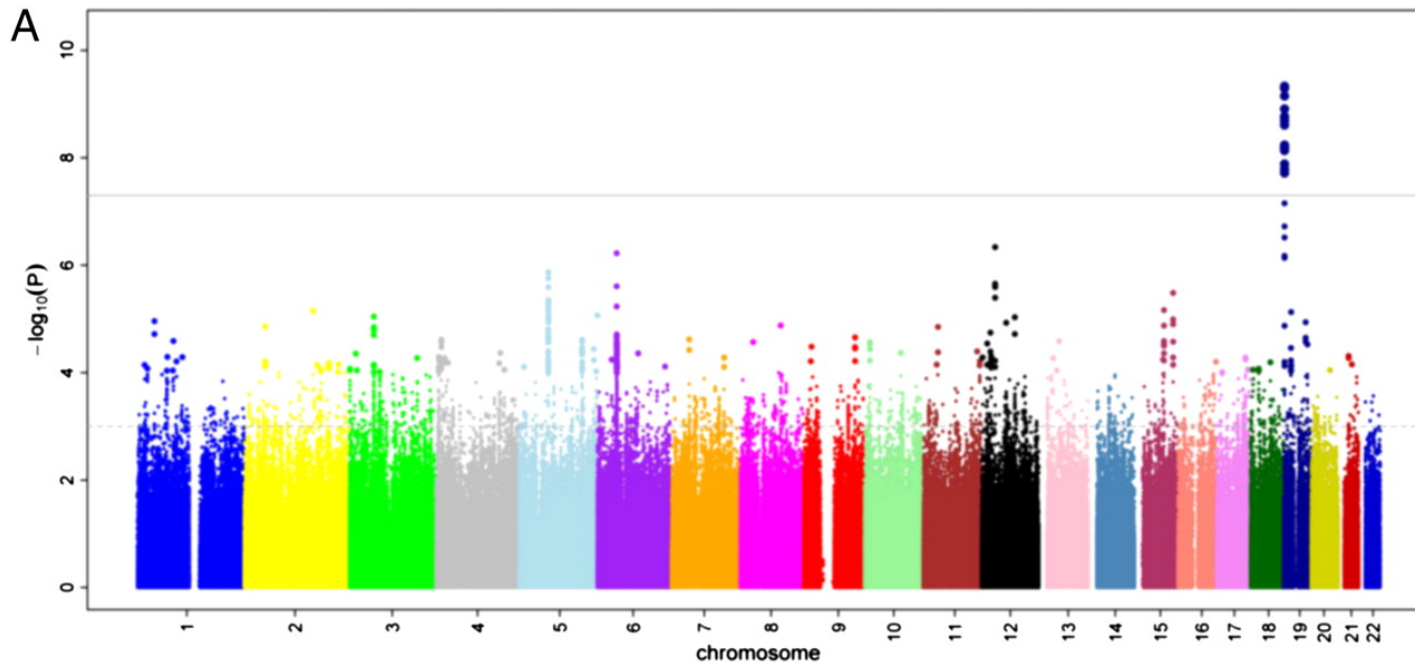


**controls (n=1,000)**  
people without heart disease



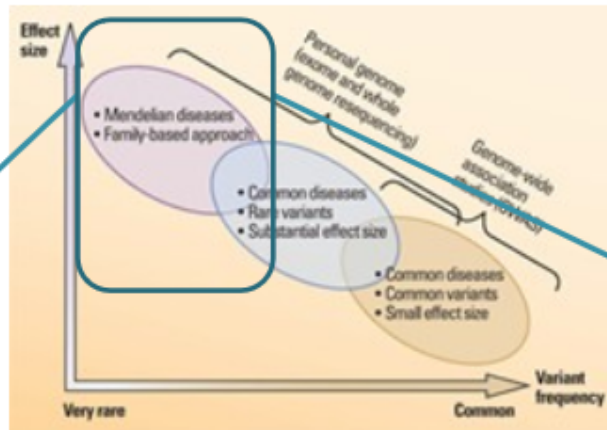
# Genome-wide association studies (GWAS)

Genome-wide association and functional studies identify the *DOT1L* gene to be involved in cartilage thickness and hip osteoarthritis





# Mendel'sche Krankheiten

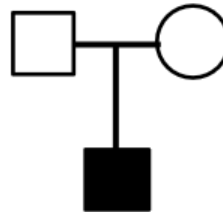


Frequency of disorder		
Rare ( $<1/10,000$ )	Low frequency ( $1/10,000$ – $1/100$ )	Common ( $>1/100$ )
Mutational target		
e.g. CHARGE syndrome ( $1/10,000$ )	e.g. Noonan syndrome ( $1/2,000$ )	e.g. intellectual disability ( $2/100$ )
<div>CHD7</div>	<div>PTPN11</div> <div>RAF1</div> <div>SOS1</div> <div>KRAS</div> <div>BRAF</div> <div>MAP2K1</div> <div>NRAS</div>	
Single gene	2–100 genes	$>100$ genes

Veltman J.A. Nat. Rev. Genetics (2012) 13:565:575.

# Mendel'sche Krankheiten: Gen-Identifizierung durch **Trio**-Analyse

- Ganze Familie wird sequenziert (gesunde Eltern sowie betroffenes Kind)



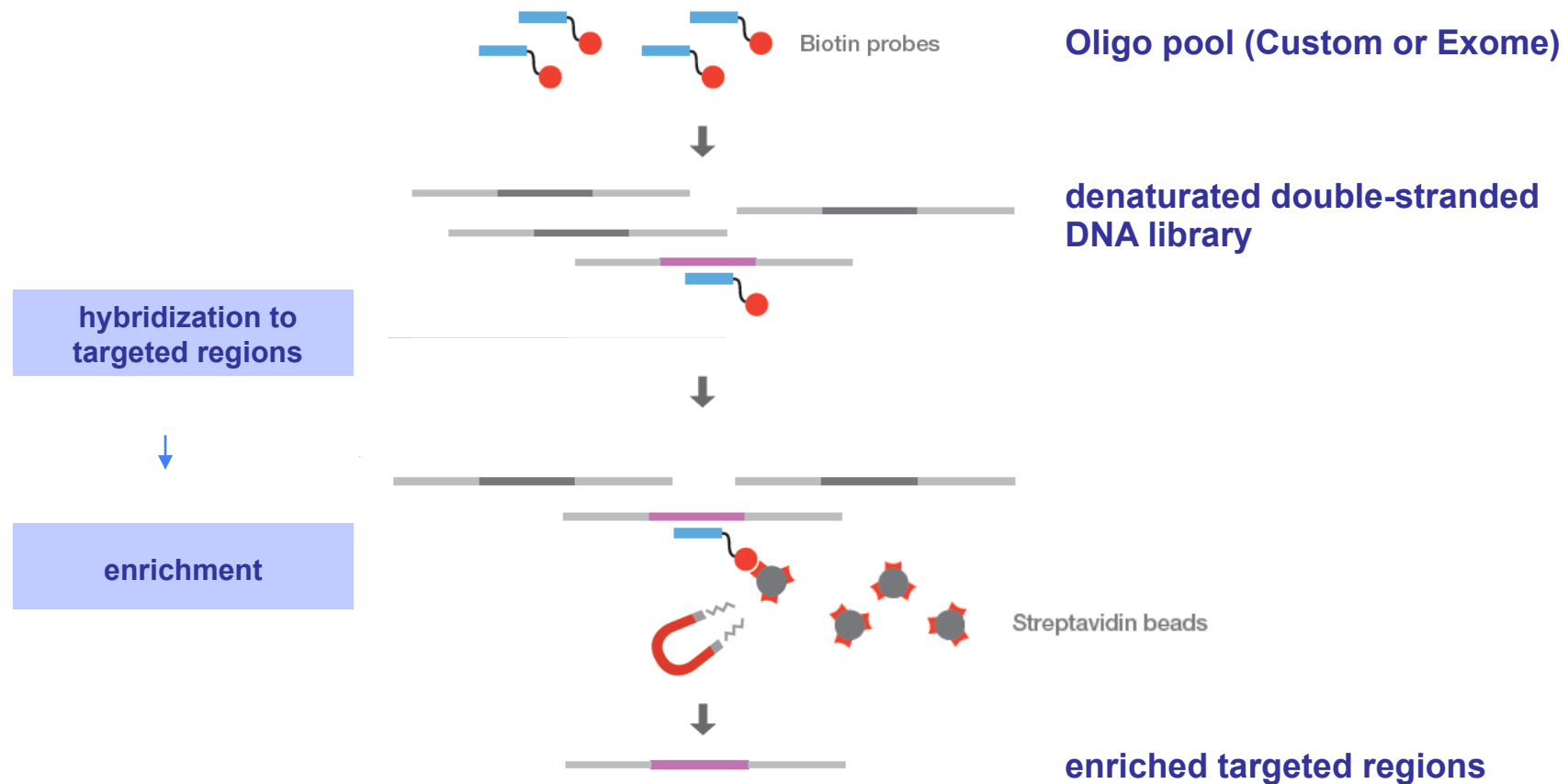
- Wieso alle drei?
  - Alle Vererbungsmodelle können untersucht werden
  - Im Gegensatz zur GWAS nur ein einziger Fall nötig!
  - Probenmaterial ist zugänglich

# Exom-Sequenzierung !!

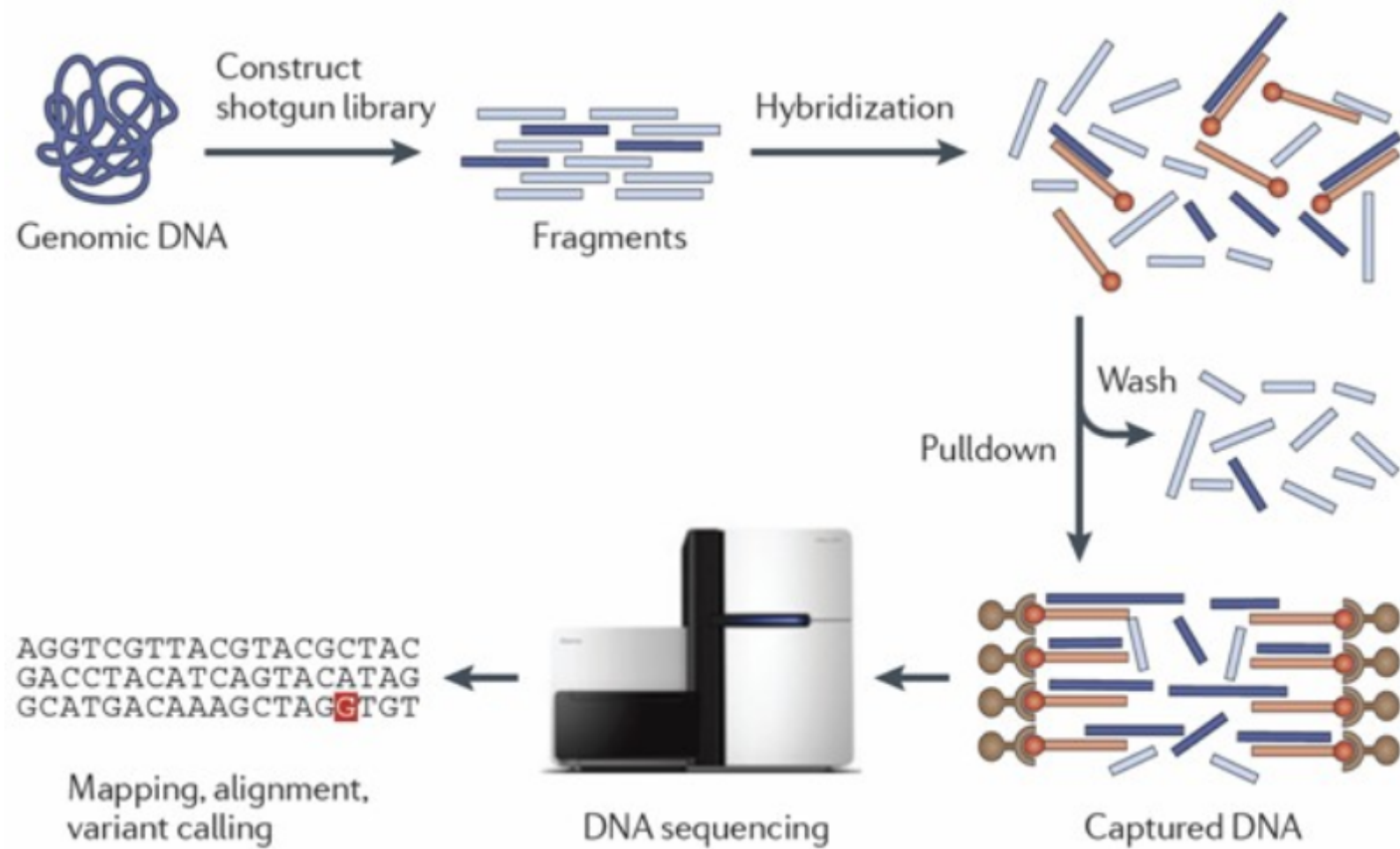
- 3.000.000.000 bp (3 Gb) Humangenom
  - Ca. 45 % repetitiv
  - Ca. 25 % Genregionen
  - Ca. 2% Exons bzw. codierende Regionen
- 20.000 – 30.000 humane Gene
  - 3000-5000 sind Krankheitsgene
  - Ca. 4000 humane genetische Krankheiten
  - 114 Gene, deren Fehlfunktion tatsächlich mit Medikamenten behandelt werden kann

# Targeted Resequencing:

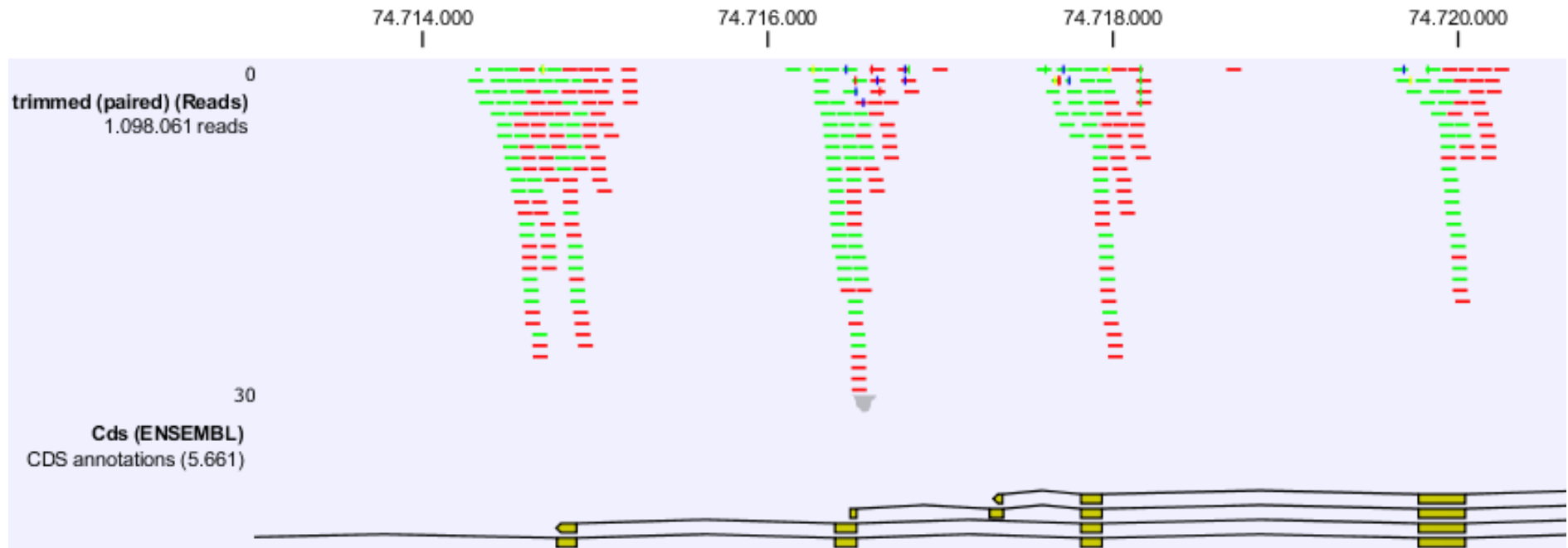
## Custom/ Exome Enrichment



# Targeted Resequencing

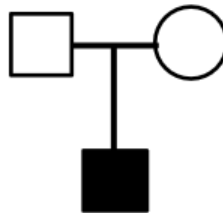


# Targeted Resequencing: Mapping



# Exomsequenzierung eines „Trios“

- Ganze Familie wird sequenziert (gesunde Eltern so wie betroffenes Kind)



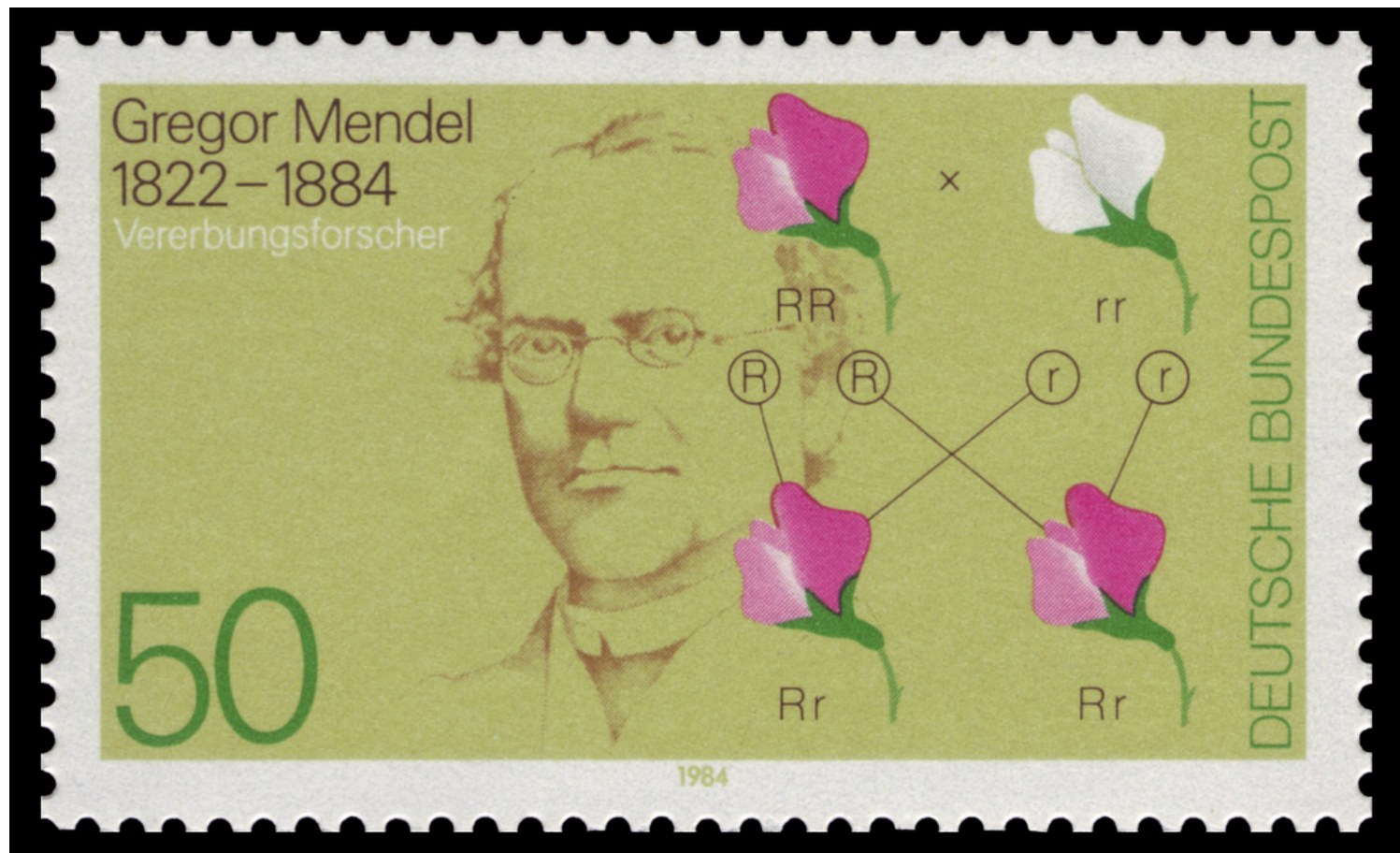
- Wieso alle 3?
  - Alle Vererbungsmodi können untersucht werden
  - Im Gegensatz zur GWAS nur ein einziger Fall nötig!
  - Probenmaterial ist zugänglich

# Die Idee dahinter?

- Identifikation von SNVs, die spezifisch für die Familie sind
- Suche nach dem „Schuldigen“ durch das Ausschlussprinzip:
  - Welche Mutationen sind neu und potenziell gefährlich?
  - An welchen Stellen wurden z.B. Varianten an das Kind homozygot weitergegeben, für die die Eltern heterozygot waren?
- Falls Parental- und Filialgeneration betroffen sind: Welche Varianten sind vom kranken Elternteil vererbt und kommen im gesunden nicht vor?
- Optimal: zusätzliche Sequenzierung von Geschwistern

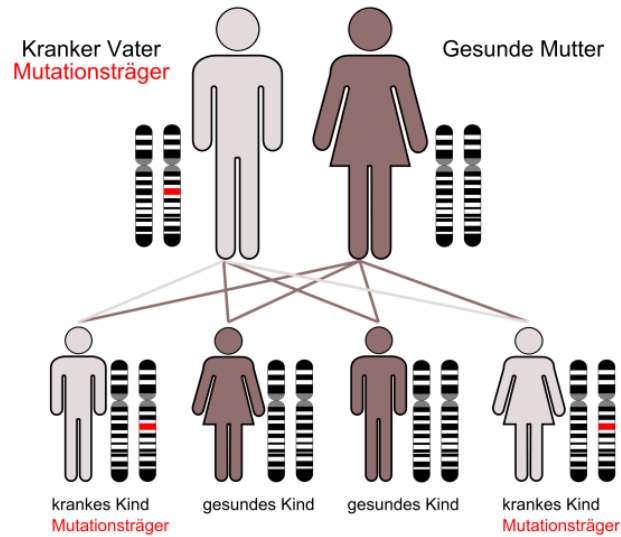


# Vererbungsmodi? Mendel!!

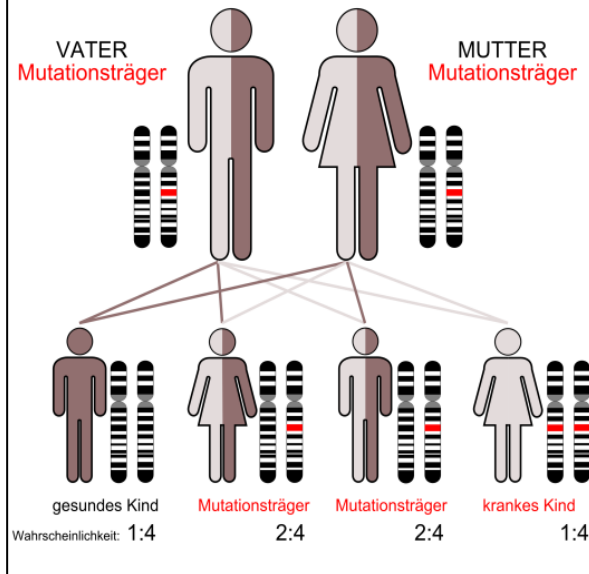


Entdecke die Möglichkeiten...

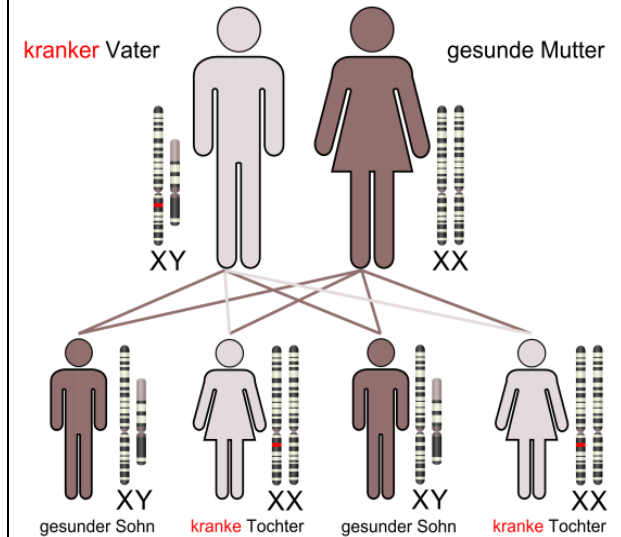
### Autosomal-dominanter Erbgang



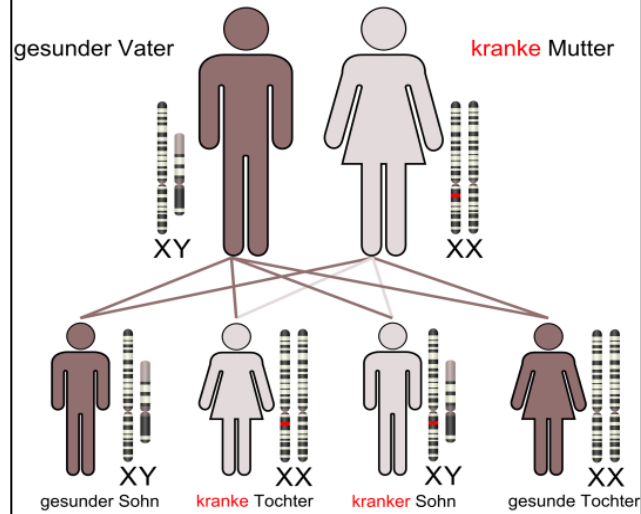
### Autosomal-rezessiver Erbgang



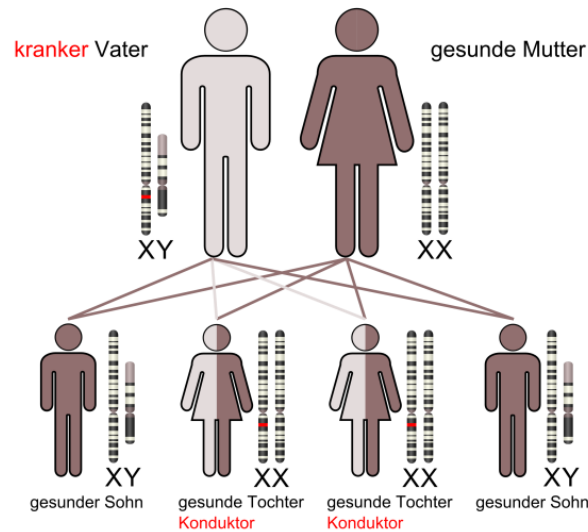
### X-chromosomal-dominanter Erbgang



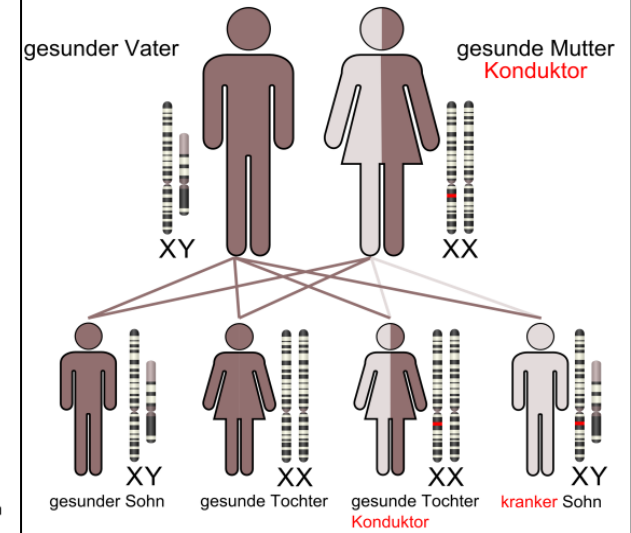
### X-chromosomal-dominanter Erbgang



### X-chromosomal-rezessiver Erbgang



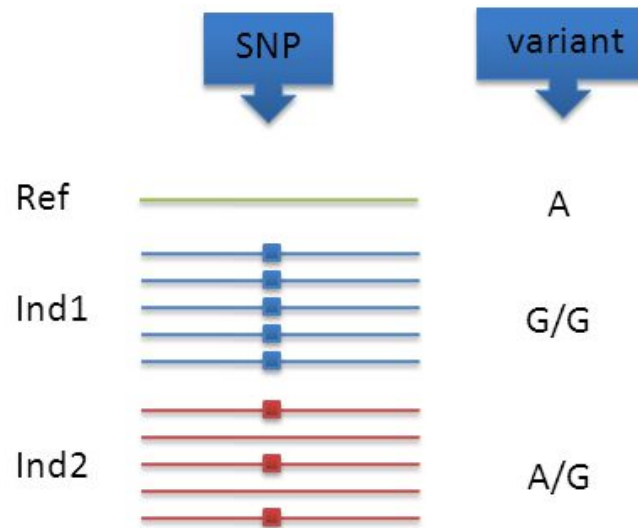
### X-chromosomal-rezessiver Erbgang



# Diagnostische Detektion von SNVs: technischer Ablauf

1) Read mapping: Kartierung der Exom-Reads gegen das humane Referenzgenom

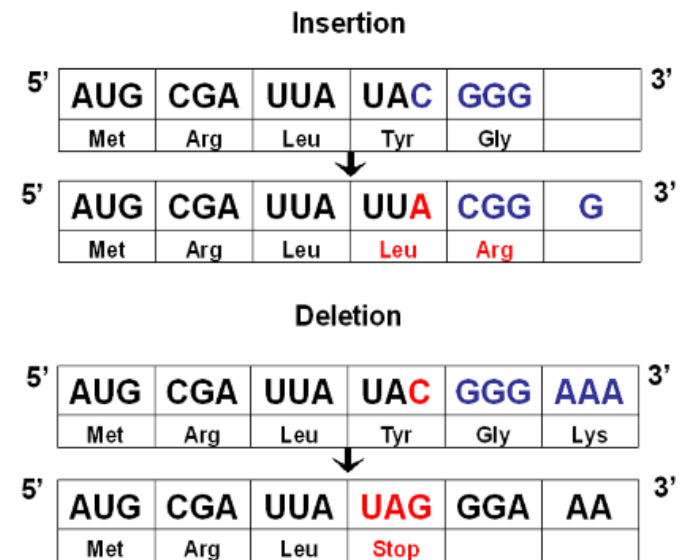
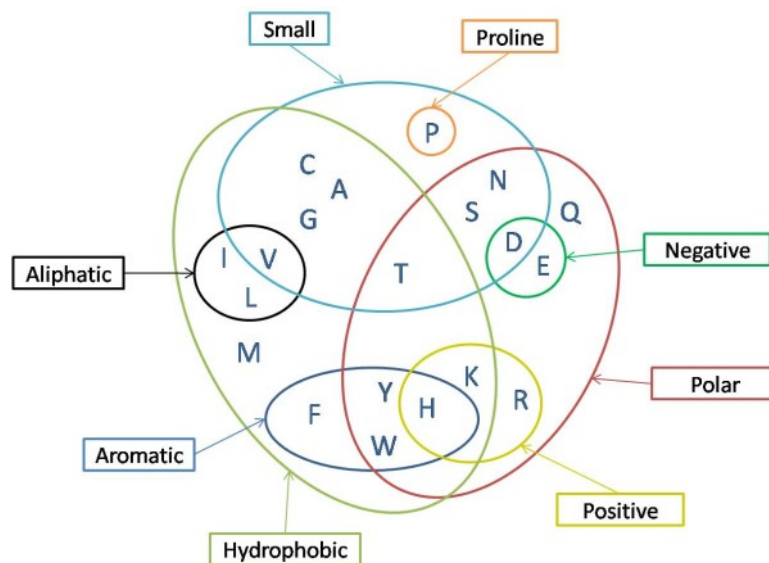
2) Identifikation der Varianten:



# Diagnostische Detektion von SNVs: technischer Ablauf

## 3) Annotation der Varianten:

- Detektion von veränderten Aminosäuresequenzen
- Abgleich mit Informationen aus öffentlichen Datenbanken

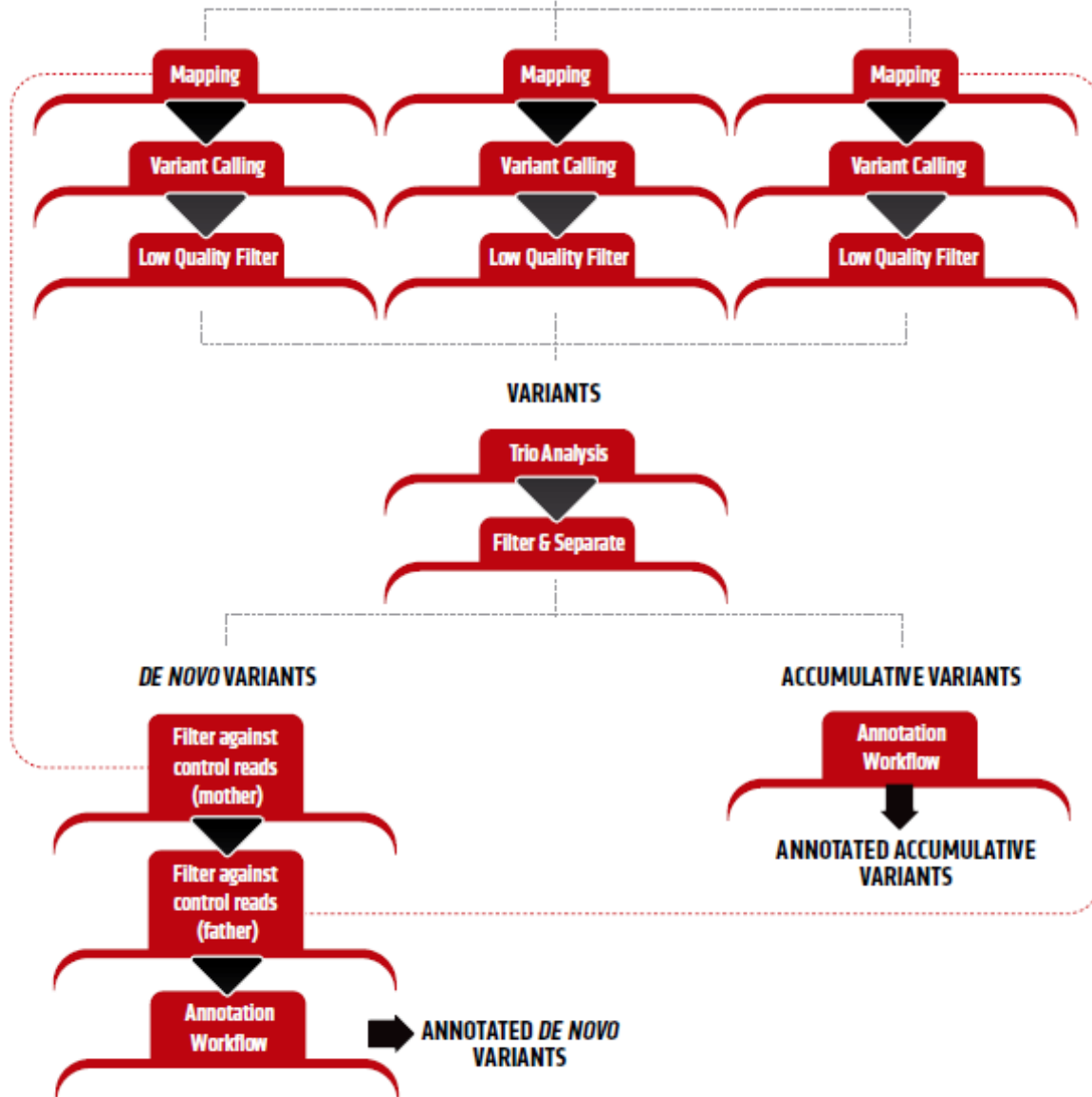




**MOTHER  
READS**

**DAUGHTER  
READS**

**FATHER  
READS**



# Technischer Ablauf

# Kursszenario

## Familie A

- 2 Kinder (2 Mädchen)
- Eines der Kinder zeigt Auffälligkeiten im Verhalten und der Entwicklung

## Familie B

- 2 Kinder (Mädchen und Junge)
- Junge zeigt Auffälligkeiten im Verhalten und der Entwicklung → große Ähnlichkeit zur ersten Familie!

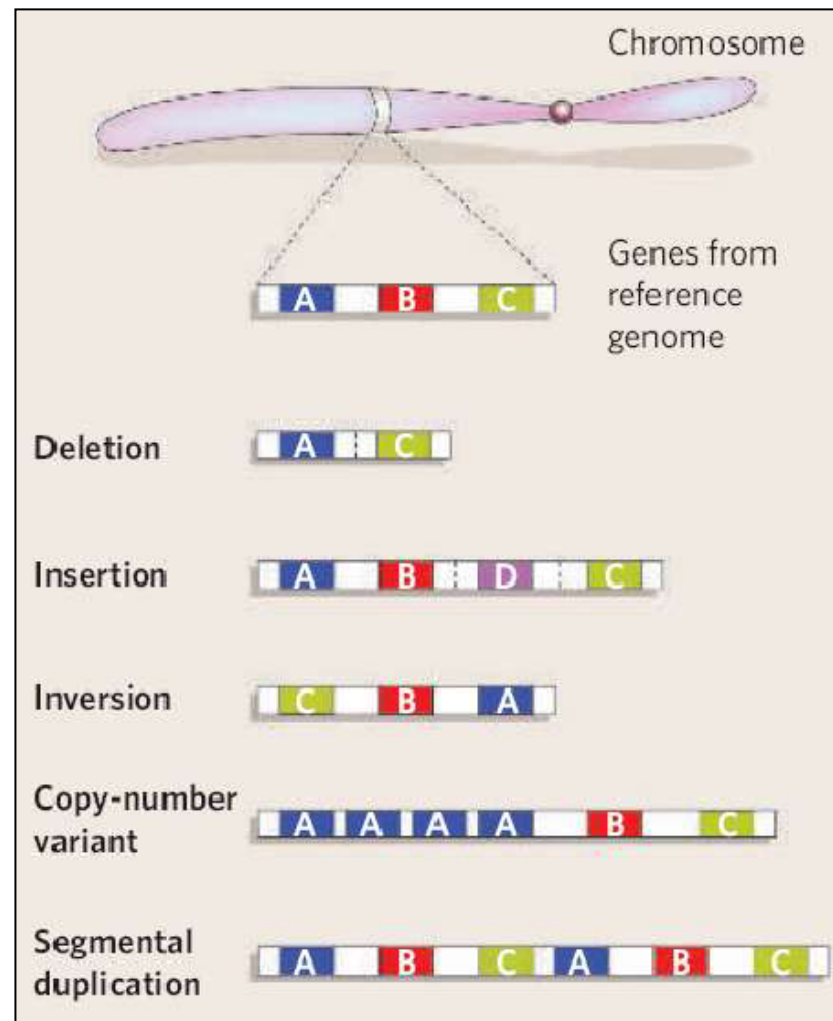
Keine Verwandtschaft zwischen den Eltern beider Familien, keine bekannte, krankheitsrelevante Vorgeschichte in beiden Fällen!

# Kursszenario

Familie 11114 (A)		Familie 11892 (B)	
<b>SRR1272259</b>	Vater	SRR1301491	Vater
<b>SRR1272260</b>	Mutter	SRR1301492	Mutter
<b>SRR1272261</b>	Krankes Kind (weibl.)	SRR1301493	Krankes Kind (männl.)
<b>SRR1272262</b>	Gesundes Kind (weibl.)	SRR1301494	Gesundes Kind (weibl.)



# Aber: SNPs/SNVs sind nicht alles!





# Copy number variations

	CNV (Database of Genomic Variants, <a href="http://projects.tcag.ca/variation/">http://projects.tcag.ca/variation/</a> )	SNP (dbSNP, <a href="http://www.ncbi.nlm.nih.gov/SNP/">http://www.ncbi.nlm.nih.gov/SNP/</a> )
Total number	38,406 <sup>a</sup> (Mar 11, 2009)	14,708,752 (Build 129)
Size	100 bp to 3 Mb	Mostly 1 bp
Type	Deletion, duplication, complex	Transition, transversion, short deletion, short insertion
Effects on genes	Gene dosage, interruption, etc.	Missense, nonsense, frameshift, splice site
Percentage of the reference genome covered	29.74% <sup>b</sup>	<1%

Mutation rates for CNVs are locus-specific and can be 100 – 10 000 fold higher than for SNPs

# Structural variation & disease

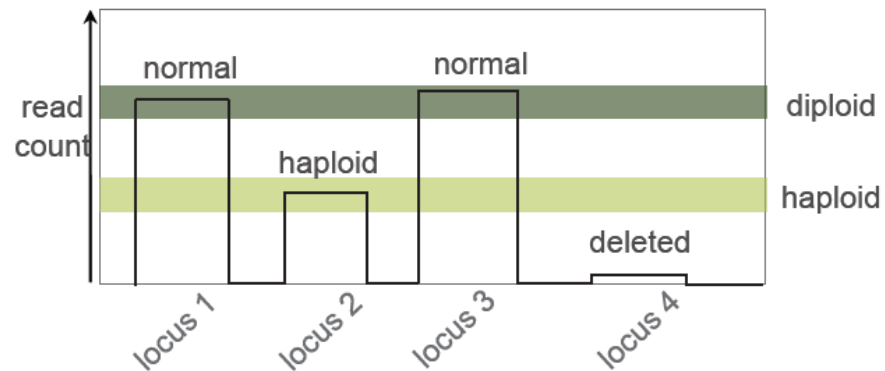
Phenotype	OMIM	Locus	CNV
<b>Mendelian (autosomal dominant)<sup>b</sup></b>			
Williams-Beuren syndrome	194050	7q11.23	del
7q11.23 duplication syndrome	609757	7q11.23	dup
Spinocerebellar ataxia type 20	608687	11q12	dup
Smith-Magenis syndrome	182290	17p11.2/ <i>RAI1</i>	del
Potocki-Lupski syndrome	610883	17p11.2	dup
HNPP	162500	17p12/ <i>PMP22</i>	del
CMT1A	118220	17p12/ <i>PMP22</i>	dup
Miller-Dieker lissencephaly syndrome	247200	17p13.3/ <i>LIS1</i>	del
Mental retardation	601545	17p13.3/ <i>LIS1</i>	dup
DGS/VCFS	188400/192430	22q11.2/ <i>TBX1</i>	del
Microduplication 22q11.2	608363	22q11.2	dup
Adult-onset leukodystrophy	169500	<i>LMNB1</i>	dup
<b>Mendelian (autosomal recessive)</b>			
Familial juvenile nephronophthisis	256100	2q13/ <i>NPHP1</i>	del
Gaucher disease	230800	1q21/ <i>GBA</i>	del
Pituitary dwarfism	262400	17q	
Spinal muscular atrophy	253300	5q1	
beta-thalassemia	141900	11p	
alpha-thalassemia	141750	16p	

Need for high-resolution detection of CNVs (and SNVs) by DNA sequencing

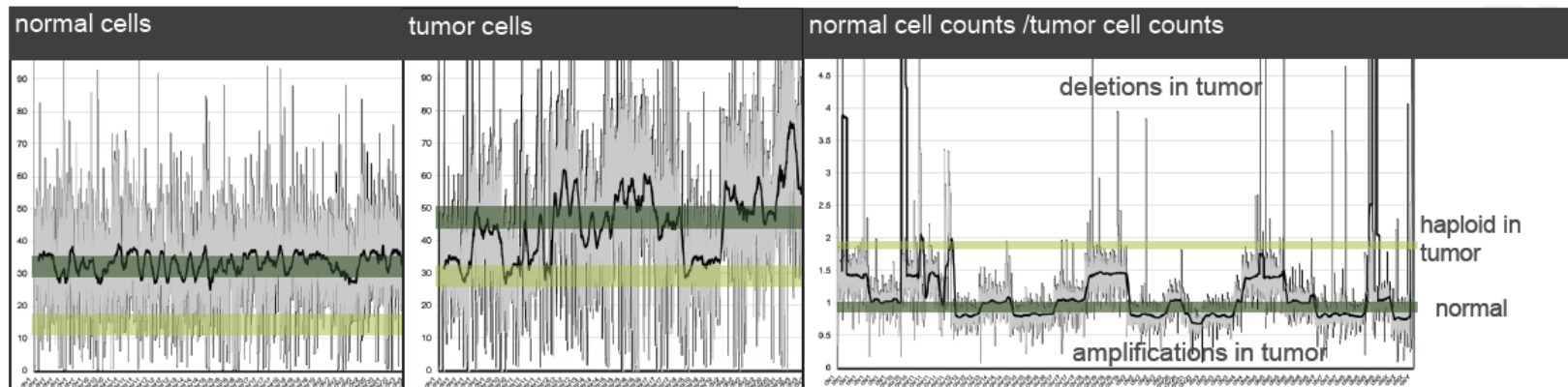
<b>Complex traits</b>			
Alzheimer disease	104300	<i>APP</i>	dup
Autism	612200	3q24	inherited homozygous del
	611913	16p11.2	del/dup
Crohn disease	266600	<i>HBD-2</i>	copy number loss
	612278	<i>IRGM</i>	del
HIV susceptibility	609423	<i>CCL3L1</i>	copy number loss
Mental retardation	612001	15q13.3	del
	610443	17q21.31	del
	300534	Xp11.22	dup
Pancreatitis	167800	<i>PRSS1</i>	tri
Parkinson disease	168600	<i>SNCA</i>	dup/tri

# Detektion von CNVs

copy number variations



Vergleich der CNV-Profile gesunder Zellen und Krebszellen eines Patienten:



# Detektion von Genom-Rearrangements

**mapping**  
of paired-end  
reads against  
genome



**filtering**  
for read pairs with  
distance outside  
expected range

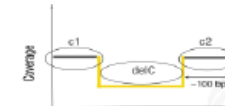


Look for  
support

Gapped  
alignment

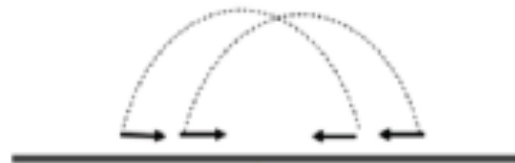


Density plots

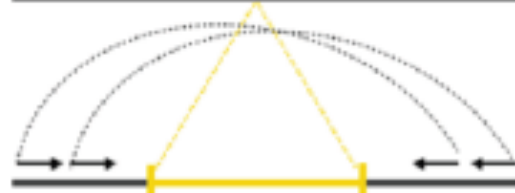


**Deletion**

Sample



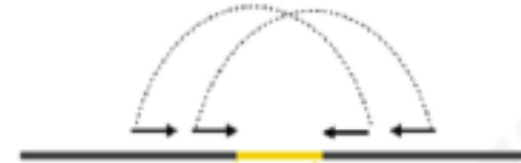
Reference



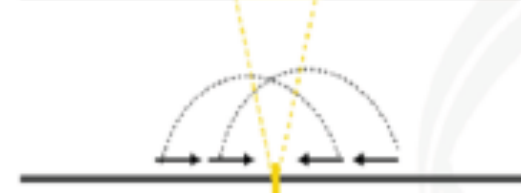
$\text{Dist} > (\mu + 3\sigma)$

**Insertion**

Sample



Reference



# Detektion von Genom-Rearrangements

**mapping**  
of paired-end  
reads against  
genome



**filtering**  
for read pairs with  
distance outside  
expected range

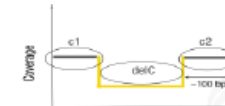


Look for  
support

Gapped  
alignment



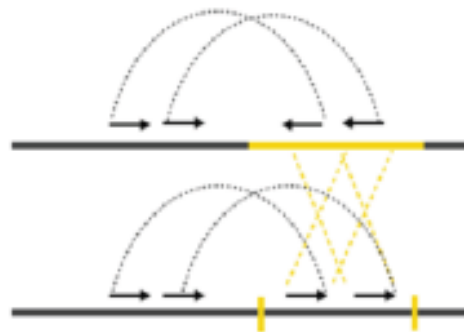
Density plots



**Inversion**

Sample

Reference



**Translocation**

ChrA

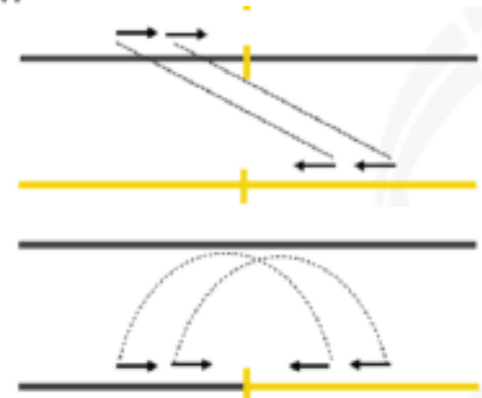
Sample

ChrB

ChrA

Reference

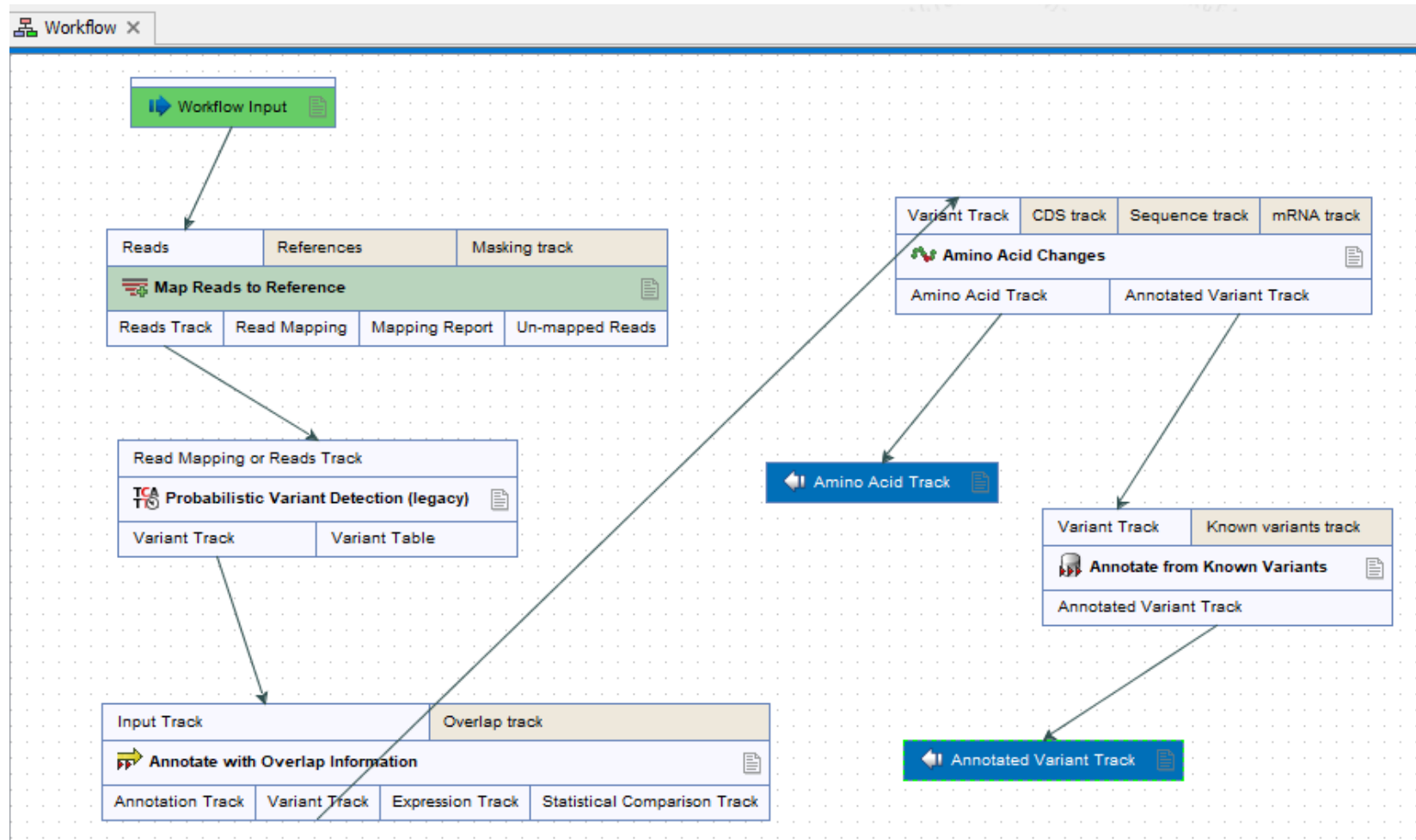
ChrB



# **Trioanalyse - Praxisteil**

## **Screenshots und Einstellungen**

# Der CLC-Workflow



# **CLC-Export:**

## **Diese Spalten müssen dabei sein!**

Position

gene\_name (Homo sapiens (hg19)\_Gene)

Amino acid change in longest transcript

Coding region change in longest transcript



Navigiere in den Ordner, in dem die Listen hinterlegt sind

Vergleiche die Geschwisterkinder untereinander und suche nach *Unterschieden*

Vergleiche die kranken Kinder von beiden Familien und finde *Gemeinsamkeiten*

```
Microsoft Windows [Version 6.3.9600]
(c) 2013 Microsoft Corporation. Alle Rechte vorbehalten.

U:\>C:
C:\>cd Users
C:\Users>cd Public
C:\Users\Public>cd "Test Triolisten"

C:\Users\Public\Test Triolisten>perl list_compare_substract_list2_from_list1 SRR1272261_denovo.csv SRR1272262_denovo.csv > Results_substr_Family1.txt

running list_compare_substract_list2_from_list1
list 1 : 'SRR1272261_denovo.csv'
list 2 : 'SRR1272262_denovo.csv'

loading file 'SRR1272262_denovo.csv' .....done!
---> loaded 740 entries

parsing file 'SRR1272261_denovo.csv' .....done!
---> parsed 421 entries
---> printed 317 entries

C:\Users\Public\Test Triolisten>perl list_compare_substract_list2_from_list1 SRR1301493_denovo.csv SRR1301494_denovo.csv > Results_substr_Family2.txt

running list_compare_substract_list2_from_list1
list 1 : 'SRR1301493_denovo.csv'
list 2 : 'SRR1301494_denovo.csv'

loading file 'SRR1301494_denovo.csv' .....done!
---> loaded 536 entries

parsing file 'SRR1301493_denovo.csv' .....done!
---> parsed 349 entries
---> printed 231 entries

C:\Users\Public\Test Triolisten>perl list_compare_find_common_entries Results_substr_Family1.txt Results_substr_Family2.txt > final_list.txt

running list_compare_find_common_entries
list 1 : 'Results_substr_Family1.txt'
list 2 : 'Results_substr_Family2.txt'

loading file 'Results_substr_Family2.txt' .....done!
---> loaded 231 entries with 209 different genes

parsing file 'Results_substr_Family1.txt' .....done!
---> parsed 317 entries
---> printed 27 entries with 20 different genes

C:\Users\Public\Test Triolisten>
```

# Detektion von Mutationen des kranken Kindes, die das gesunde Geschwisterkind nicht hat

```
1  #!/usr/bin/perl
2
3  ##### check input #####
4
5  $0 =~ s/.*\.(.+)/$1/;
6  unless (@ARGV >= 2)
7  -{
8      print STDERR "\n\n$0\n";
9      print STDERR "- x length ($0);
10     print STDERR "\nscript to compare two tab separated lists, subtract list2 & prints entries from list 1 when conform with column 1 & 2\n";
11     print STDERR "results given on STDOUT (>)\n";
12     print STDERR "\n";
13     print STDERR "USAGE\n";
14     print STDERR "$0 [list1] [list2]\n";
15     print STDERR "\n\n";
16     exit;
17 }
18
19 my $list1 = $ARGV[0];
20 my $list2 = $ARGV[1];
21
22 die "\n$list1: no such file!\n\n" unless (-f $list1);
23 die "\n$list2: no such file!\n\n" unless (-f $list2);
24
25 ##### /checkinput #####
26
27 print STDERR "\nrunning $0\n";
28 print STDERR "list 1 : '$list1'\n";
29 print STDERR "list 2 : '$list2'\n";
30
31 my %LIST1;
32 my $list1_entry_sum = 0;
33 my $list1_printed = 0;
34 my %LIST2;
35 my $list2_entry_sum = 0;
36
37 print STDERR "\nloading file '$list2'...";
38
39 open (DAT, "<", "$list2") || die "\ncouldn't open file '$list2'!\n\n";
40
41 while (<DAT>)
42 -{
43     chomp ($_);
44     my @L = split (/t/, $_);
45
46     ##### WORK IN HERE #####
47
48     $LIST2{$L[1]}{$L[0]}++;
49     $list2_entry_sum++;
50
51     ##### /WORK IN HERE #####
52 }
53
54 close (DAT);
55
56 print STDERR "\rloading file '$list2'.....done!\n";
57 print STDERR "----> loaded $list2_entry_sum entries\n";
58
59 print STDERR "\nparsing file '$list1'...";
60
61 my %GENES_PRINTED_SUM;
62
63 open (DAT, "<", "$list1") || die "\ncouldn't open file '$list1'!\n\n";
64
65 while (<DAT>)
66 -{
67     chomp ($_);
68     my @L = split (/t/, $_);
69
70     ##### WORK IN HERE #####
71
72     $list1_entry_sum++;
73
74     next if (exists ($LIST2{$L[1]}) && exists ($LIST2{$L[1]}{$L[0]}));
75
76     print "$_ \n";
77     $list1_printed++;
78     $GENES_PRINTED_SUM{$L[1]}++;
79
80     ##### /WORK IN HERE #####
81 }
82
83 close (DAT);
84
85 print STDERR "\rparsing file '$list1'.....done!\n";
86 print STDERR "----> parsed $list1_entry_sum entries\n";
87 print STDERR "----> printed $list1_printed entries\n\n";
88
```

# Schnittmenge der mutierten Gene beider kranken Kinder

```
1  #!/usr/bin/perl
2
3  ##### checkinput #####
4
5  $0 =~ s/.*\.(.+)/$1/;
6  unless (@ARGV >= 2)
7  {
8      print STDERR "\n\n$0\n";
9      print STDERR "-" x length ($0);
10     print STDERR "\nscript to compare two tab separated lists and print entries from list 1 sharing common column (2) with list 2\n";
11     print STDERR "results given on STDOUT (>)\n";
12     print STDERR "\n";
13     print STDERR "USAGE\n";
14     print STDERR "$0 {list1} {list2}\n";
15     print STDERR "\n\n";
16     exit;
17 }
18
19 my $list1 = $ARGV[0];
20 my $list2 = $ARGV[1];
21
22 die "\n$list1: no such file!\n\n" unless (-f $list1);
23 die "\n$list2: no such file!\n\n" unless (-f $list2);
24
25 ##### /checkinput #####
26
27 print STDERR "\nrunning $0\n";
28 print STDERR "list 1 : '$list1'\n";
29 print STDERR "list 2 : '$list2'\n";
30
31 my %LIST1;
32 my $list1_entry_sum = 0;
33 my $list1_printed = 0;
34 my %LIST2;
35 my $list2_entry_sum = 0;
36
37 print STDERR "\nloading file '$list2'...";
38
39 open (DAT, "<", "$list2") || die "\ncouldn't open file '$list2'!\n\n";
40
41 while (<DAT>)
42 {
43     chomp ($_);
44     my @L = split (/[, \_]/, $_);
45
46     ##### WORK IN HERE #####
47
48     $LIST2{$L[1]}{$L[2]}++;
49     $list2_entry_sum++;
50 }
```

```
51     ##### /WORK IN HERE #####
52 }
53
54 close (DAT);
55
56 print STDERR "\rloading file '$list2'.....done!\n";
57 print STDERR "---> loaded $list2_entry_sum entries with ".keys (%LIST2)." different genes\n";
58
59 print STDERR "\nparsing file '$list1'...";
60
61 my %GENES_PRINTED_SUM;
62
63 open (DAT, "<", "$list1") || die "\ncouldn't open file '$list1'!\n\n";
64
65 while (<DAT>)
66 {
67     chomp ($_);
68     my @L = split (/[, \_]/, $_);
69
70     ##### WORK IN HERE #####
71
72     $list1_entry_sum++;
73
74     if (exists ($LIST2{$L[1]}))
75     {
76         print "$_ \n";
77         $list1_printed++;
78         $GENES_PRINTED_SUM{$L[1]}++;
79     }
80
81     ##### /WORK IN HERE #####
82 }
83
84 close (DAT);
85
86 print STDERR "\rparsing file '$list1'.....done!\n";
87 print STDERR "---> parsed $list1_entry_sum entries\n";
88 print STDERR "---> printed $list1_printed entries with ".keys (%GENES_PRINTED_SUM)." different genes\n\n";
```