

# **Exome sequencing at the CCG**

***“NGS an der Schnittstelle von  
Grundlagen- und translationaler  
Forschung”***



**Peter Nürnberg**

**TMF-Workshop  
Berlin, 7.12.2012**

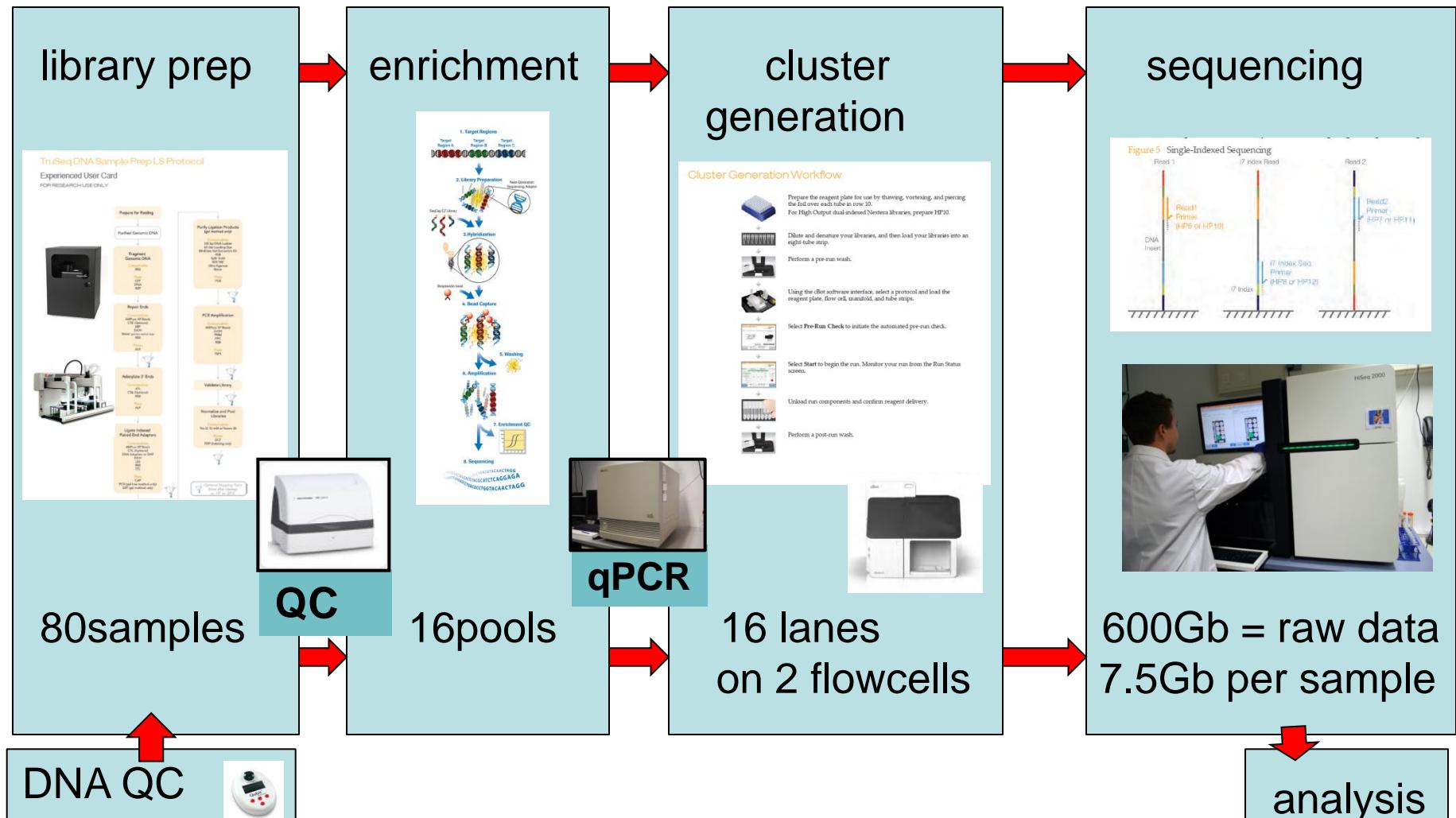


Cologne Center for  
Genomics

University of Cologne



# The exome sequencing pipeline



Cologne Center for  
Genomics



# From DNA to library

- DNA QC and input in the data base
- Library prep  
(fragmentation, end repair , adenylate 3`ends, Adapter ligation, purification, amplification, validation)



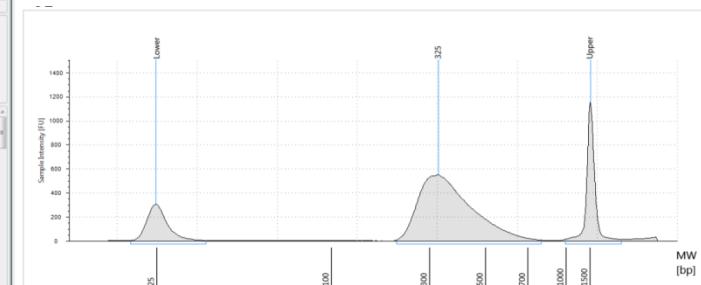
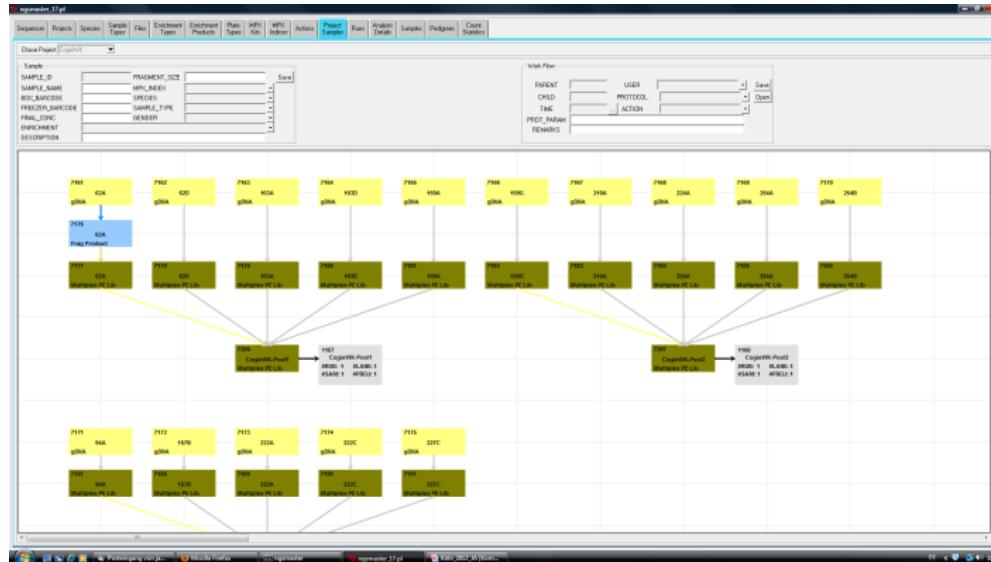
Beckman Coulter. Biomek® FXP

illumina®

TruSeq® DNA  
Sample Preparation Guide



FOR RESEARCH USE ONLY  
ILLUMINA PROPRIETARY  
Part # 15026480 Rev. C  
July 2012



G01\_Lib

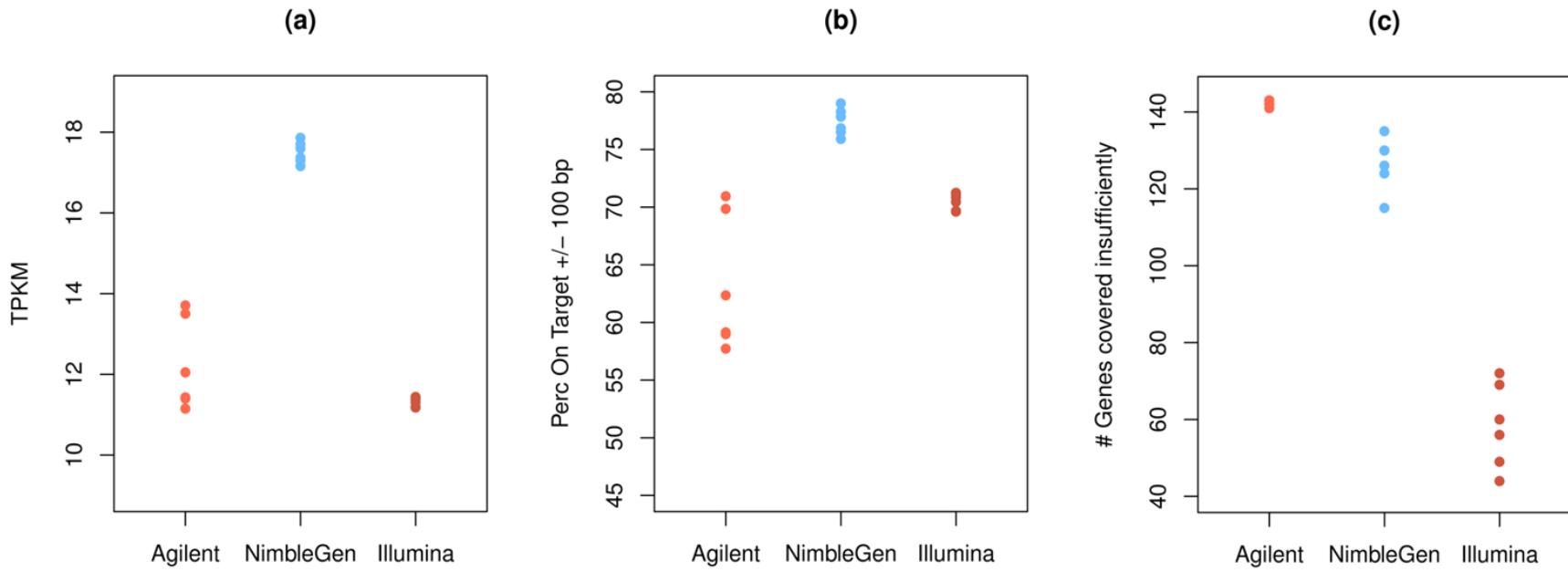
Wavelength	MW [bp]	Conc. (ng/μl)	Molarity (nmol/l)	% of Integrated Area	Peak Comment	Observations
Sample	25	4.46	271			Lower Marker
Sample	325	20.8	96.9	100.00		
Sample	1,500	6.50	6.57			Upper Marker

# *Enrichment: different products*

company	products	Size and content	Ease of use	Performance Q1/2011	Price Q1/2011
Agilent	V1	38Mb	x		
	50Mb	51Mb	x(+)	(3.)	
	V4	51Mb	x		
	V4+UTR	71Mb	x		
Nimblegen	Array	34Mb	x		
	V2	44Mb	x	(1.)	(+)
	V3	64Mb	x		
Illumina	TruSeq Nextera	62Mb	x	(2.)	

# *Assessing the enrichment performance in targeted resequencing experiments*

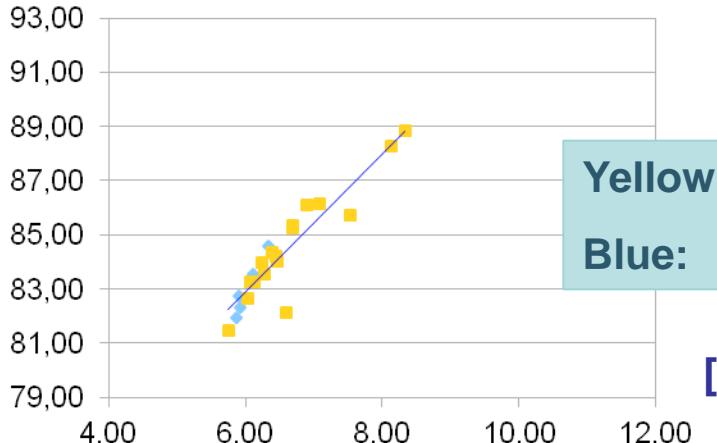
Peter Frommolt et al., Hum Mut 2011



TPKM = On-**T**arget reads **P**er **K**ilobase target region and **M**illion mappable reads

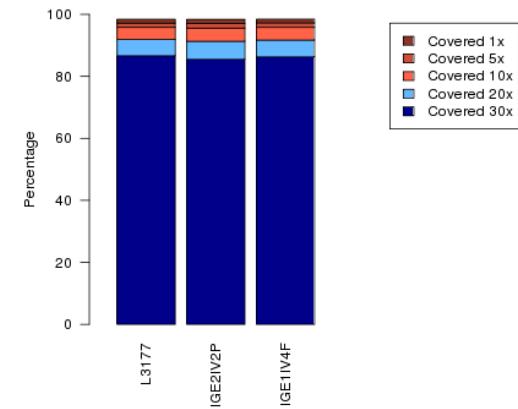
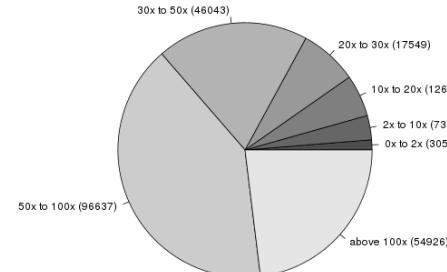
# Nimblegen v2: strategy considerations

% of target > 30X



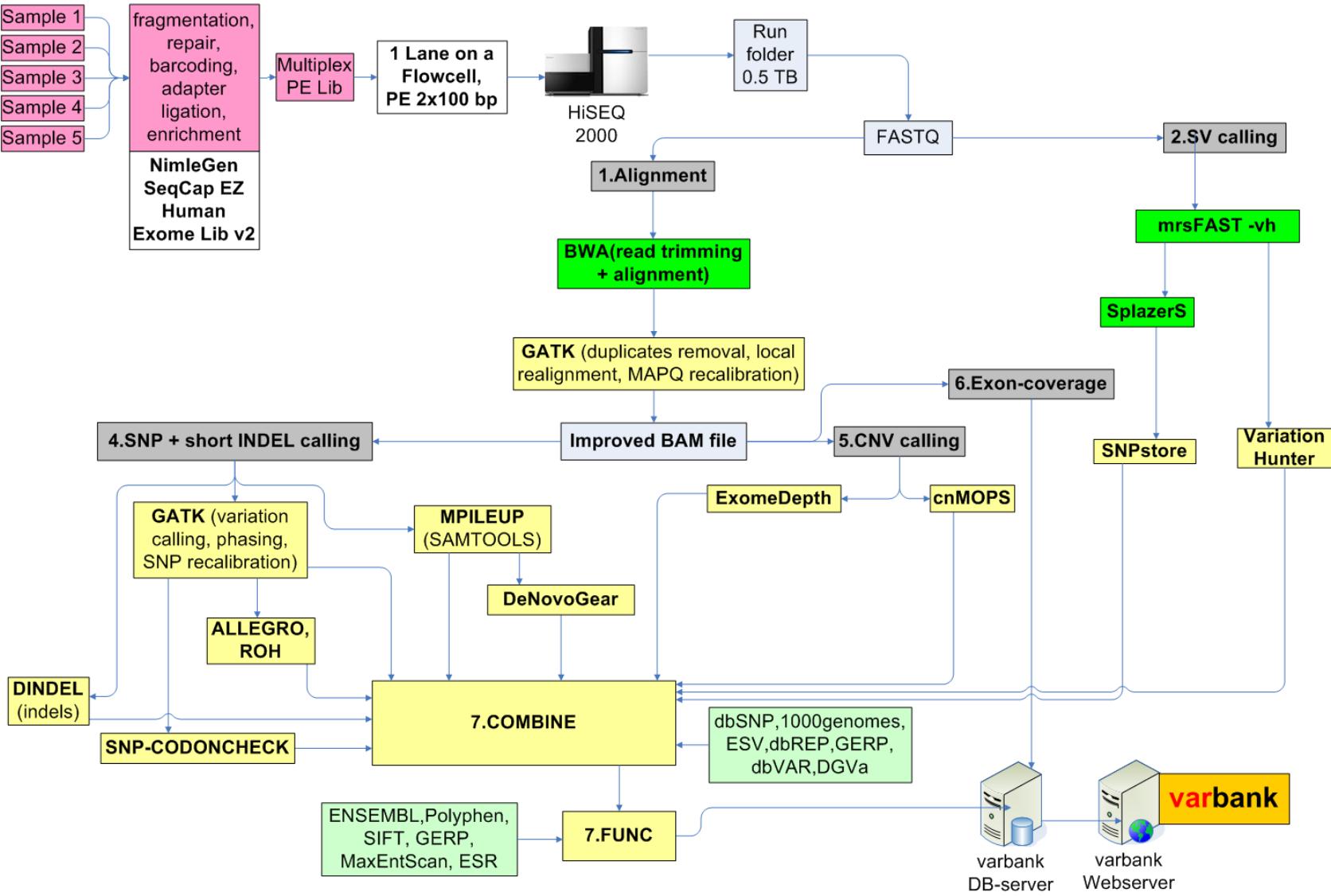
Yellow: pre-capture pooling  
Blue: post-capture pooling

	ideal	real
Raw data	7.5Gb	5-10Gb
aligning	(all)	90-92%
Aligning to target sequence	(all)	72-78%



# Reads	64620968
# On Target ± 100 bp	44549948
Target Size (bp)	44234111
# Target Regions	238236
Coverage Mean	77.27
Coverage Std Dev	64.9
Covered 1x	97.93%
Covered 5x	96.46%
Covered 10x	94.61%
Covered 20x	89.55%
Covered 30x	82.45%
TPKM	15.59

# Exome Sequencing Pipeline at the Cologne Center for Genomics



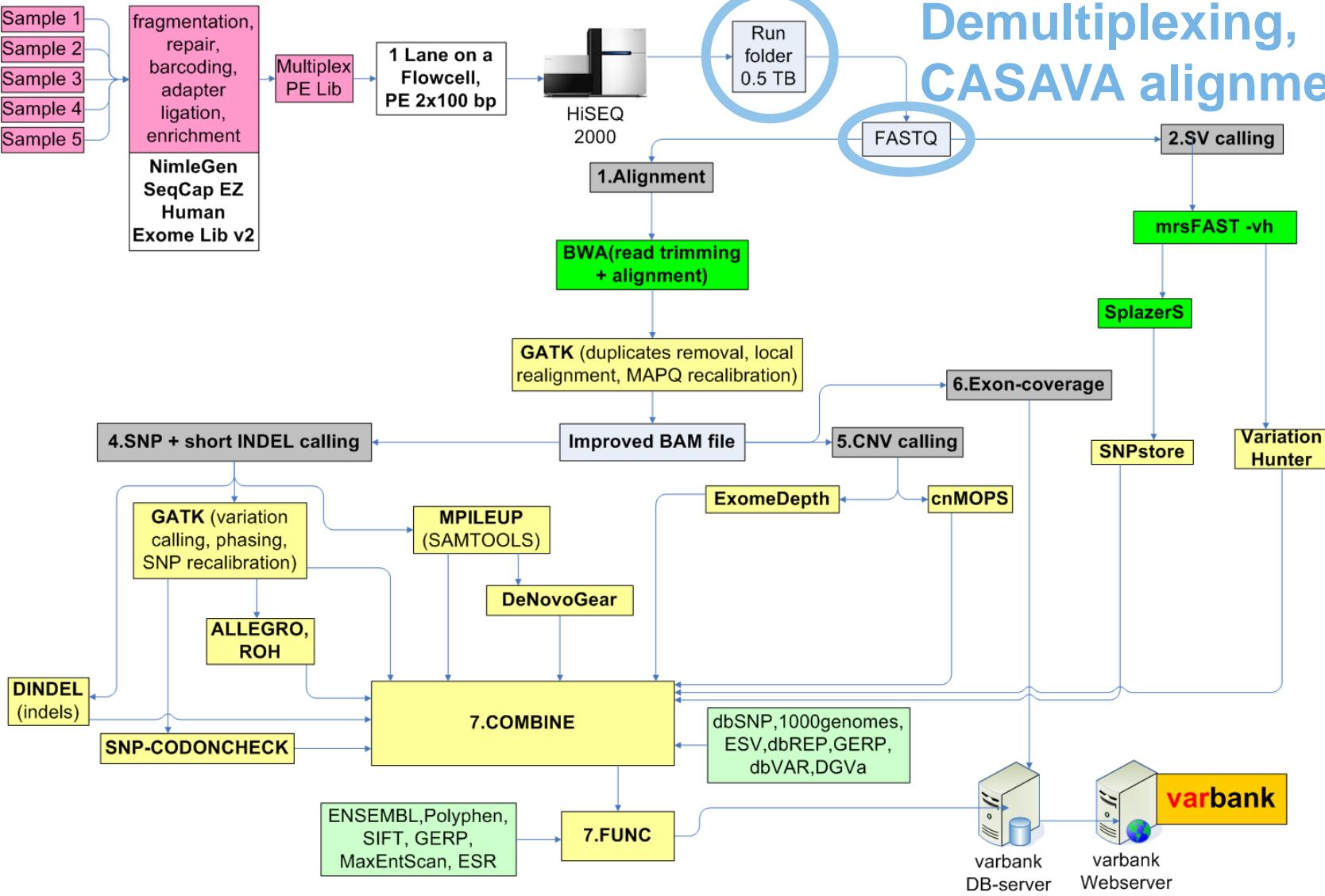
# **Goals for a new exome analysis pipeline**

- Establishing analysis pipeline on HPC Clusters available at the Regional Computing Center Cologne (CHEOPS/SUGI)
- Improving alignments
- Expanding the range of SNP/short indel callers
- Adding support for structural variations and CNVs
- Building up a new database and web tools to have remote access to variants, files, statistics, etc.
- Design the web tool to make variant browsing, filtering and interpretation as easy as possible

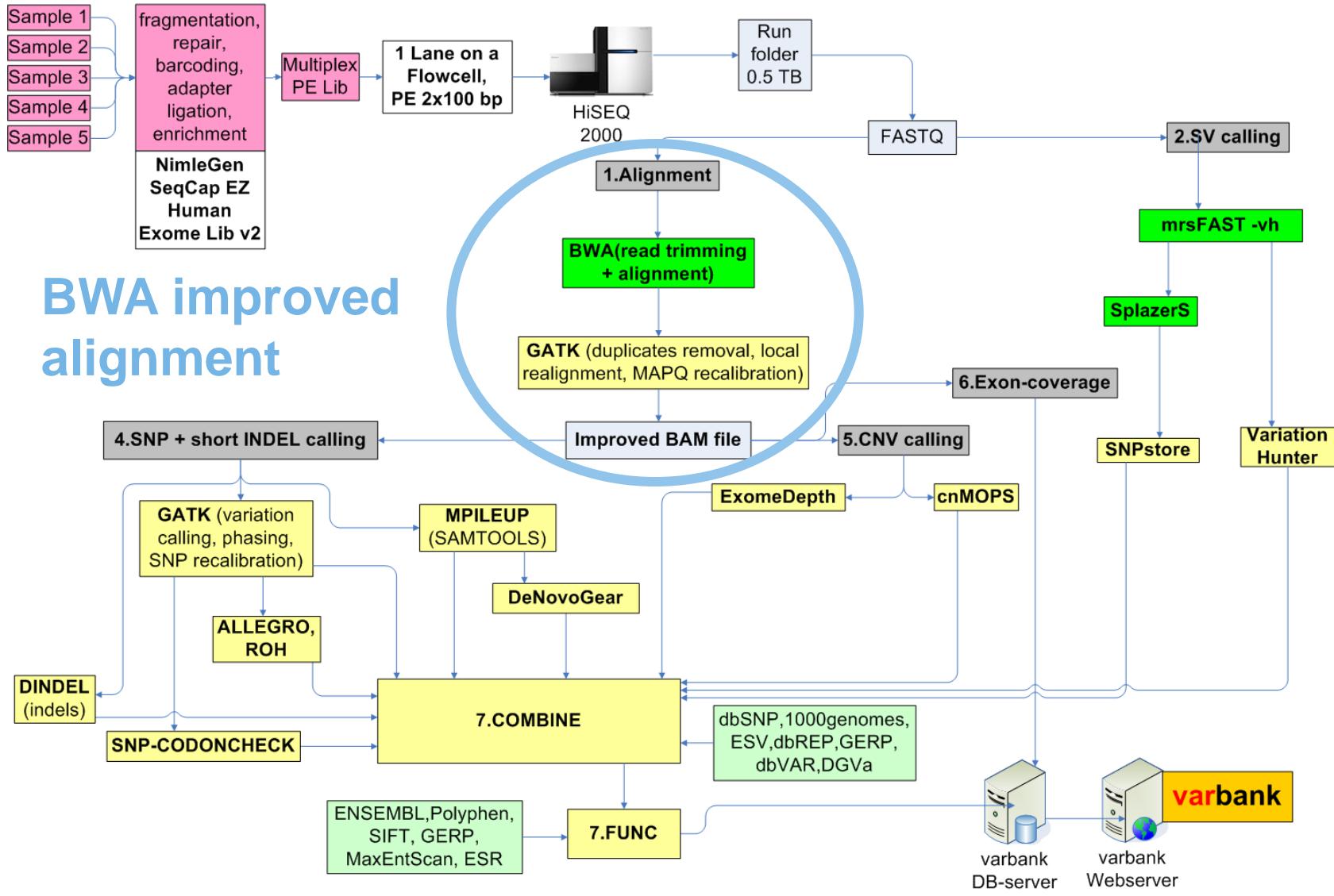


# Exome Sequencing Pipeline at the Cologne Center for Genomics

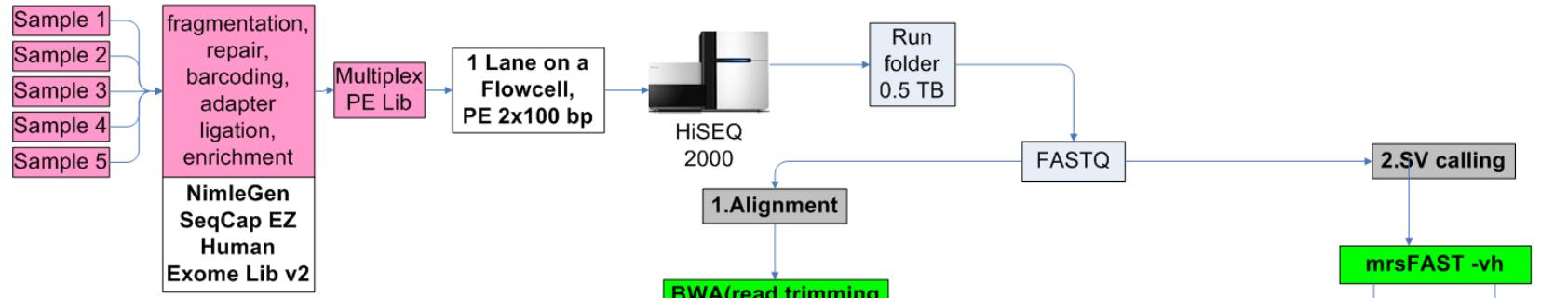
## Demultiplexing, CASAVA alignment



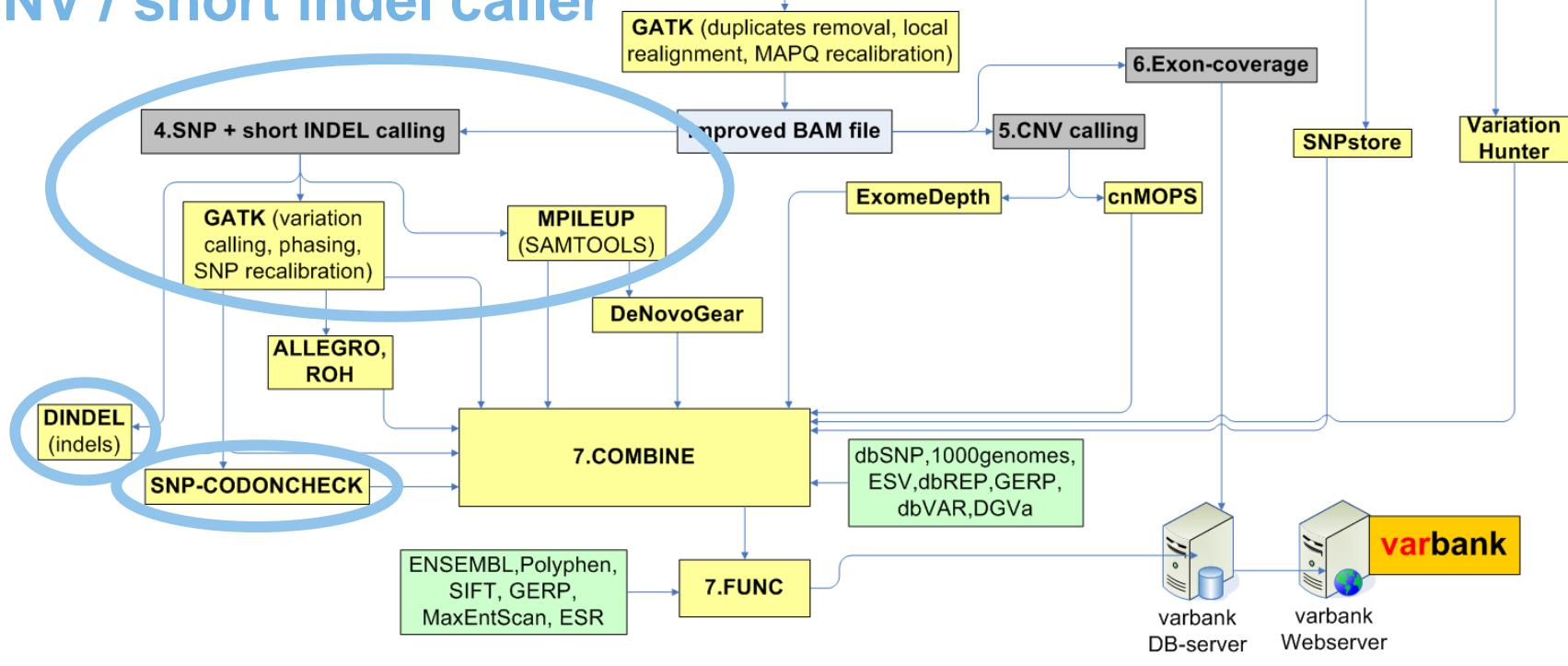
# Exome Sequencing Pipeline at the Cologne Center for Genomics



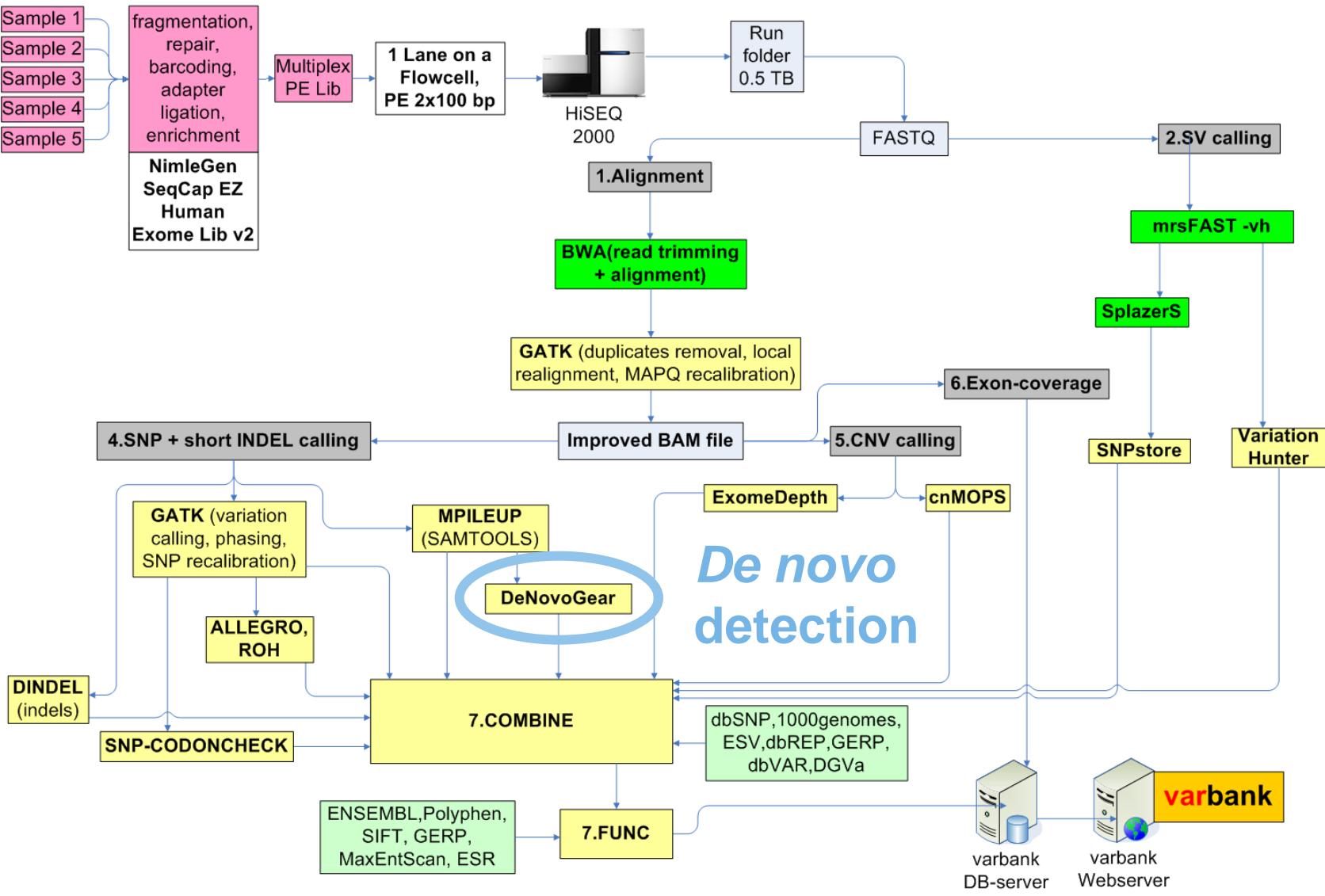
# Exome Sequencing Pipeline at the Cologne Center for Genomics



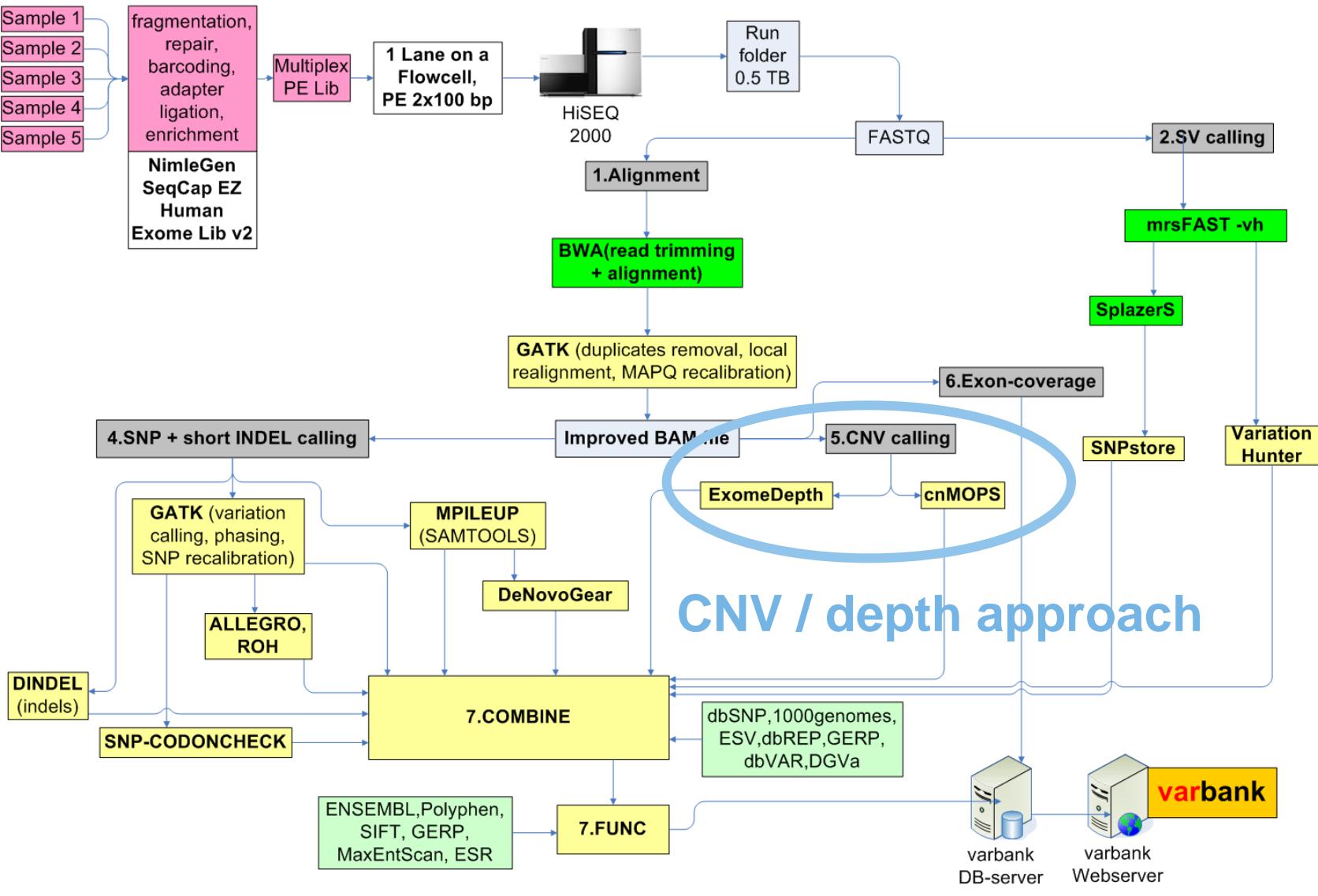
## SNV / short indel caller



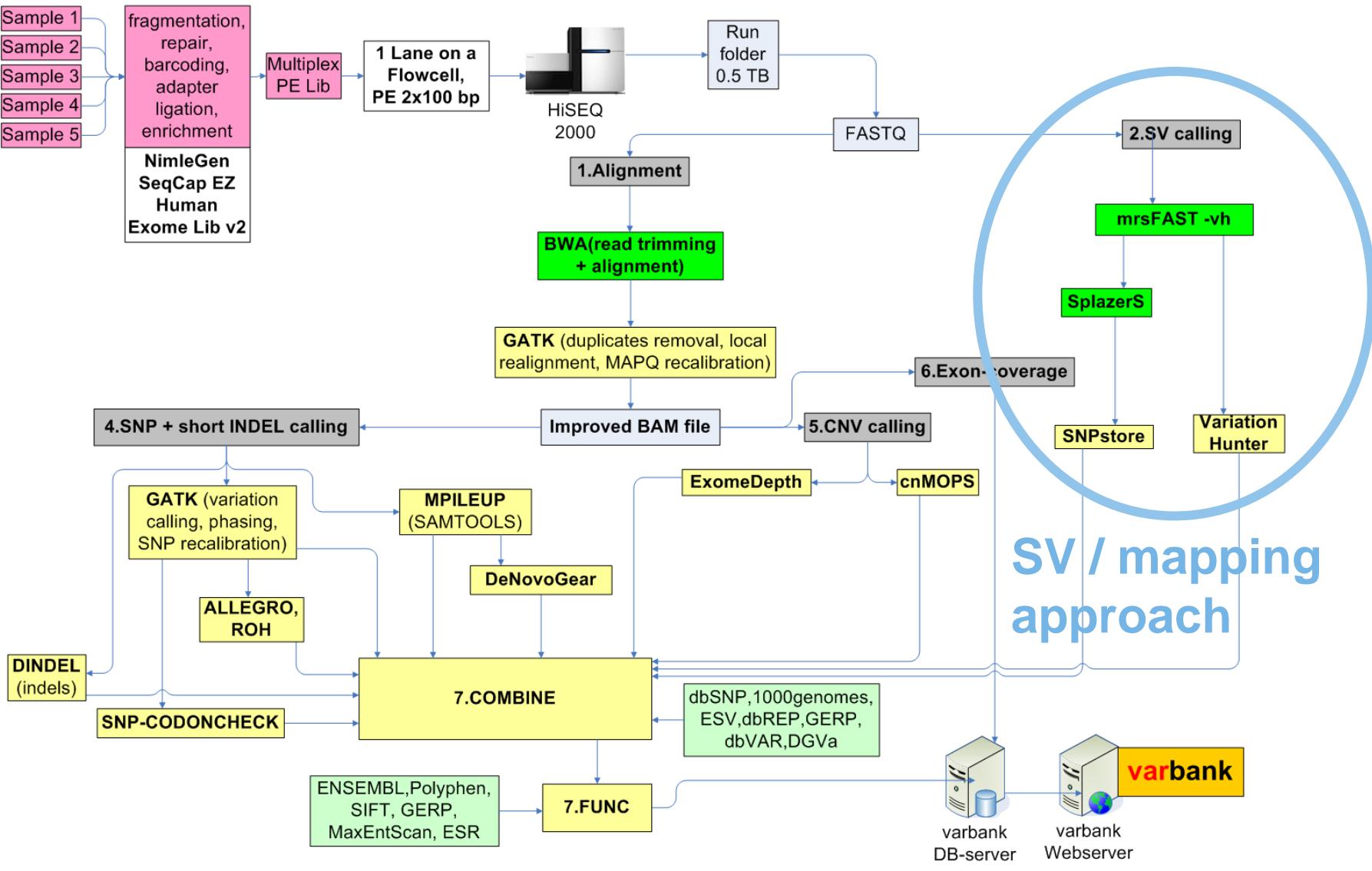
# Exome Sequencing Pipeline at the Cologne Center for Genomics



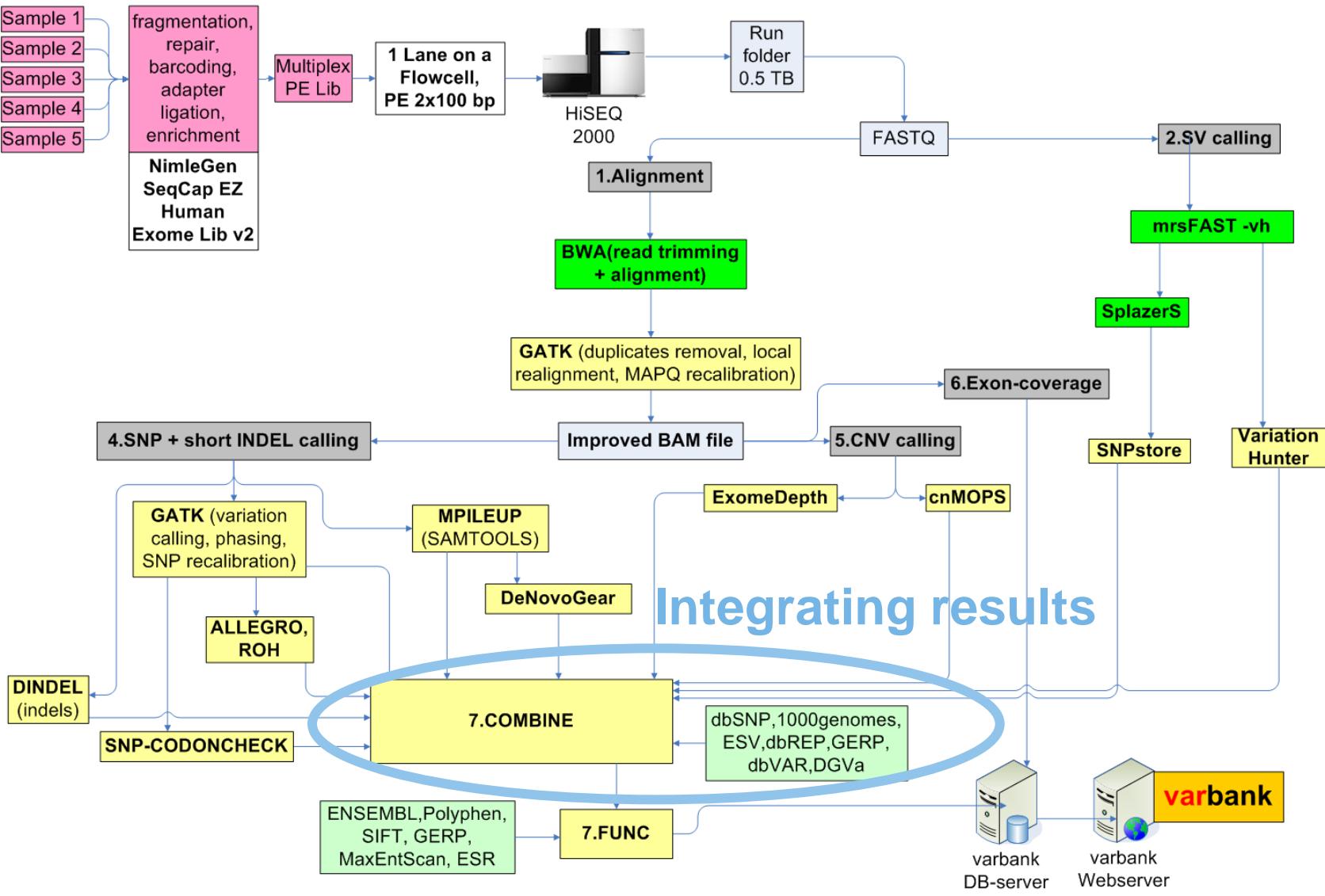
# Exome Sequencing Pipeline at the Cologne Center for Genomics



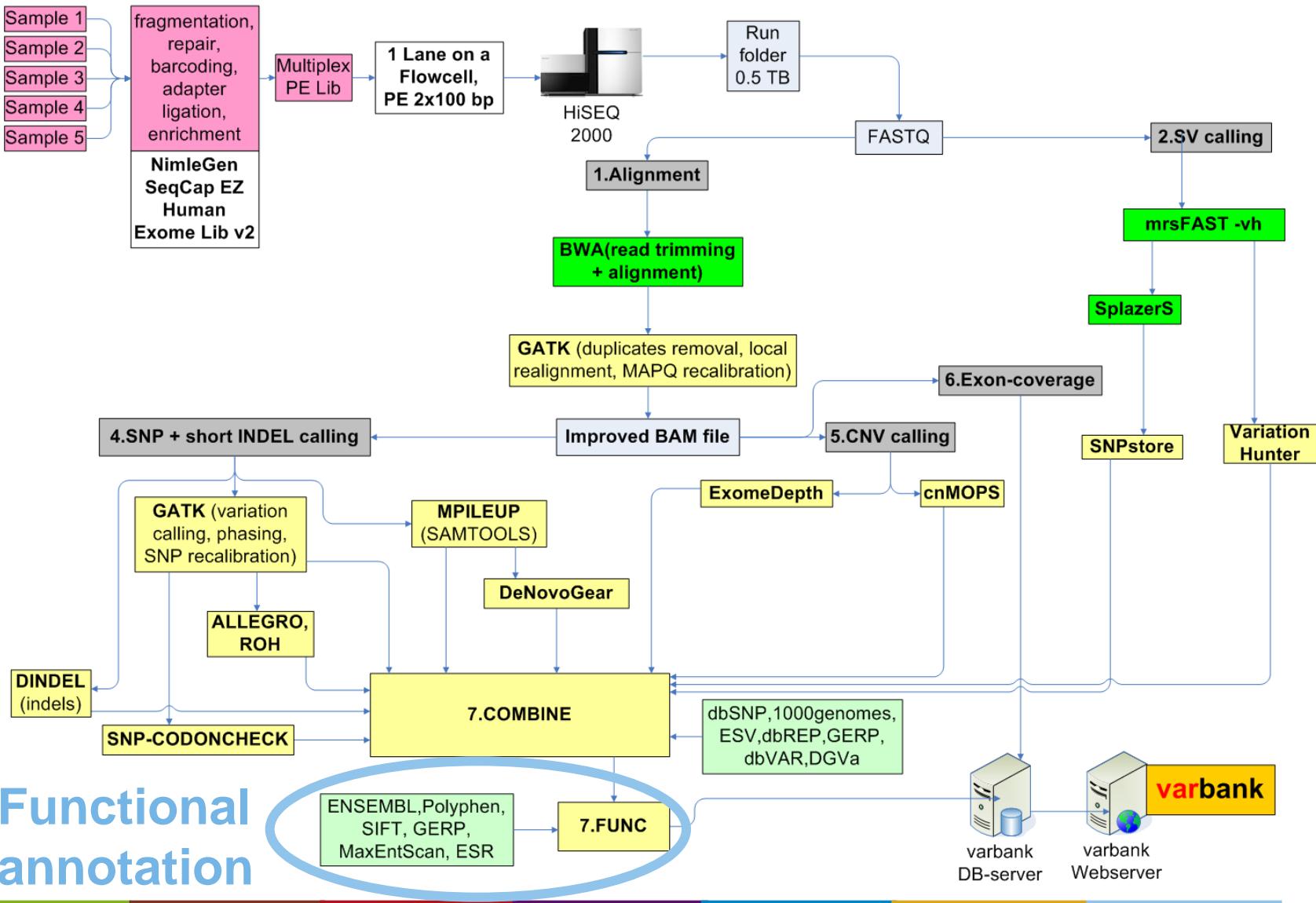
# Exome Sequencing Pipeline at the Cologne Center for Genomics



# Exome Sequencing Pipeline at the Cologne Center for Genomics



# Exome Sequencing Pipeline at the Cologne Center for Genomics



Functional  
annotation

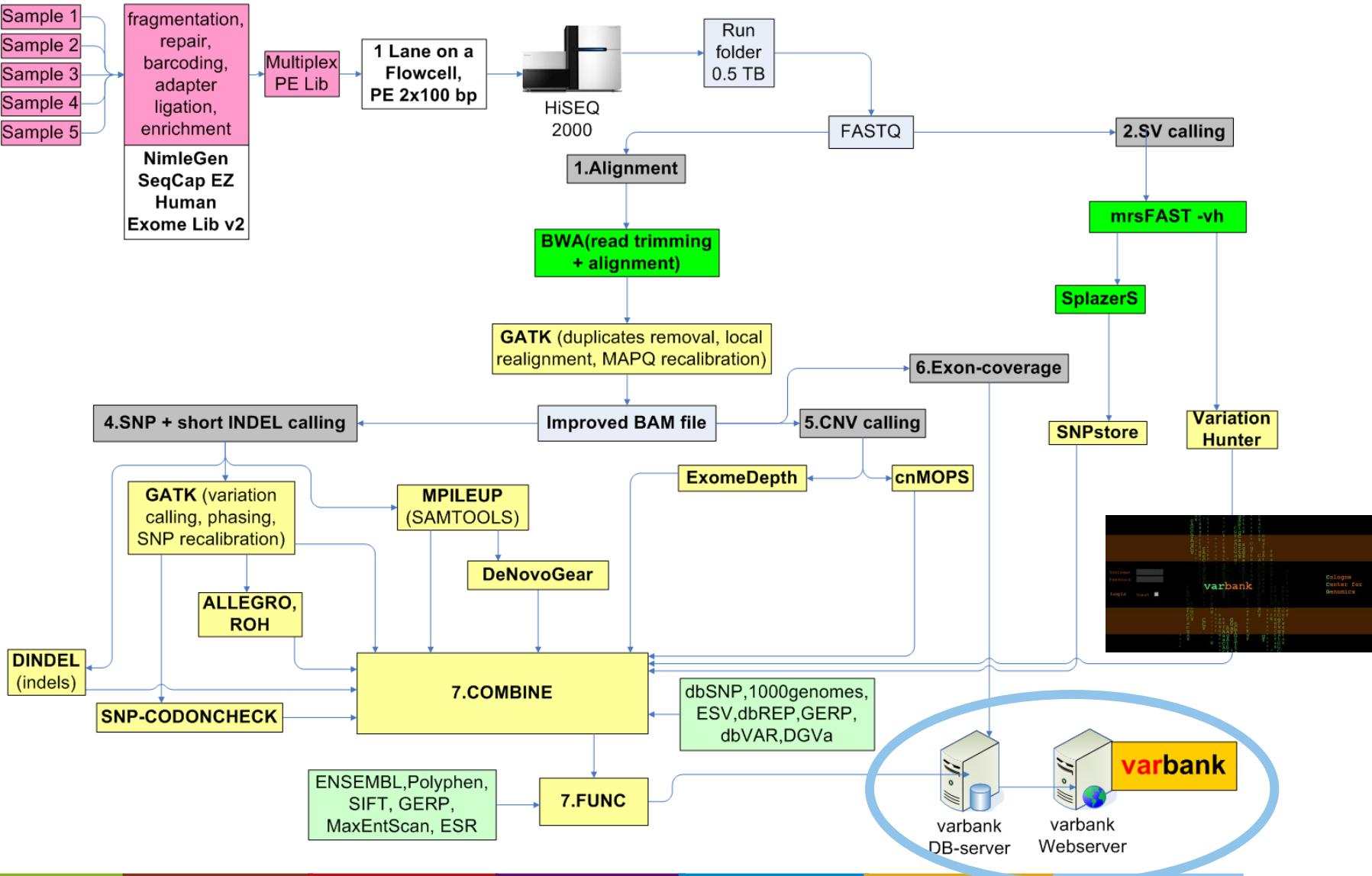


Cologne Center for  
Genomics

Universität zu Köln



# Exome Sequencing Pipeline at the Cologne Center for Genomics



Username   
Password

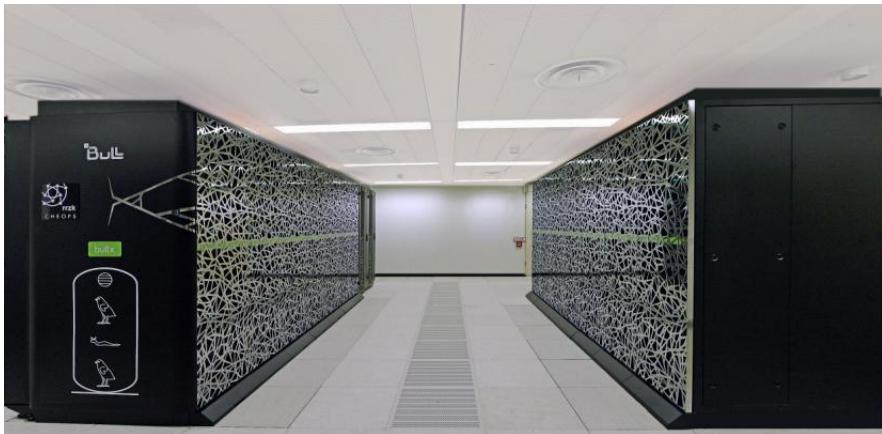
Login Guest



varbank

Cologne  
Center for  
Genomics

# CHEOPS: HPC@cologne



- 100 Tflop/s : in production since 2010
- 210 x 2 Nehalem EP quad-core processors 24 GB RAM
- 432 x 2 Westmere hexa-core processors 24 GB RAM
- 170 x 2 Westmere hexa-core processors 48 GB RAM
- 5 x 2 Nehalem EP quad-core processors 96 GB RAM
- 24 x 4 Nehalem EX octo-core processors 512 GB RAM
- Lustre file system, 500 TB

# *varbank - features*

- **Repository for all relevant data files**
- **Fine-tunable access control (login, sample owner, sample sharing)**
- **User-editable for information related to pedigrees and phenotypes**
- **Variant filtering, browsing and exporting**
- **Basic read alignment viewer**
- **Gene coverage viewer**

# **varbank – technical implementation**

- Implementation on two VMWare servers located at the Regional Computing Center Cologne (outsourced maintenance, scalable!)
- Data base: Oracle 11.2 with table space partitioning
- Web server: Apache
- Web interface: perl/cgi; java script; jQuery; jqgrid; jQuery.svg; overlibmws; Matrix.js; ~ 5200 code lines

	AB	Prid	Aid	Sid	Enrich	Pipeline	GenBuild	EnsBuild	Pedigree	Sample	Project	Owner	Phenotype	SubPhenTypes
	A	265	1206	5274	Nimbl-SeqCapEZ-V2 37M	2.1_devel	hg19	B68		ROL_0111	RW	Fritz Zimprich	Rolando	2
	A	265	1201	5265	Nimbl-SeqCapEZ-V2 37M	2.1_devel	hg19	B68		ROL_0091	RW	Fritz Zimprich	Rolando	2
	A	265	1200	5253	Nimbl-SeqCapEZ-V2 37M	2.1_devel	hg19	B68		ROL_0041	RW	Fritz Zimprich	Rolando	2
	A	265	1211	5729	Nimbl-SeqCapEZ-V2 37M	2.1	hg19	B68		ROL_0671	RW	Fritz Zimprich	Rolando	2

Home Dataset Selection/Filter Options Browse Exomes Browse Reads Gene Coverage Download Exon Coverage Documentation

Variation Data

	EnsProt	Strand	Biotype	HGNC	CCDS	RefSeq	Dist5SS	Dist3SS	RegulationInfo	MutPos	PostFilter	Consequence	Con
A 26													
A 24													
A 24	LZ ENSP00000382526	1	protein_coding	CACNA1C	CCDS53733.1	NM_001167624..-89	-15	ESR=20,6,8,13,\$3=	CDS.43	CDS	COMPLEX_INDEL;ESE2ESS_CHANGE	X--X-	
34	ENSP00000382542	1	protein_coding	CACNA1C		NM_001167625..1389	584	BPO=CCGAT,2,3,3,INT.43	INTRON	BP_LOSS;CRYPTIC_5SS_ACTIVATION;E-X-XX			
12	ENSP00000329877	1	protein_coding	CACNA1C	CCDS53736.1	NM_001129830..-89	-15	ESR=20,6,8,13,\$3=	CDS.43	CDS	COMPLEX_INDEL;ESE2ESS_CHANGE	X--X-	
52	ENSP00000353868	-1	protein_coding	PCED1A					INT.3	INTRON	DEC_3SS_STRONG;ESE_LOSS	-X-X-	
32	ENSP00000353072	1	protein_coding	GNRH2									
42	ENSP00000369705	1	protein_coding	GNRH2									
33	ENSP00000245983	1	protein_coding	GNRH2									
46	ENSP00000369704	1	protein_coding	GNRH2									
00	ENSP00000352003	1	protein_coding	GNRH2									
47	ENSP00000369705	1	protein_coding	GNRH2									
32	ENSP00000353072	1	protein_coding	GNRH2									
00	ENSP00000352003	1	protein_coding	GNRH2									
33	ENSP00000245983	1	protein_coding	GNRH2									
46	ENSP00000369704	1	protein_coding	GNRH2									
72	ENSP00000347184	1	protein_coding	HTT									
34	ENSP00000448683	1	protein_coding	IL32									
13	ENSP00000371648	1	protein_coding	IL32									
30	ENSP00000433747	1	protein_coding	IL32									
58	ENSP00000324742	1	protein_coding	IL32									
15	ENSP00000405063	1	protein_coding	IL32									
34	ENSP00000450364	1	protein_coding	IL32									
30	ENSP00000008180	1	protein_coding	IL32									
52	ENSP00000446624	1	protein_coding	IL32	CCDS32379.1	NM_001012636..-308	ESR=44,15,44,23,\$	CDS.6					
33	ENSP00000411958	1	protein_coding	IL32	CCDS32377.1	NM_004221.4	-308	ESR=44,15,44,23,\$	CDS.7				

[Wikipedia](#) [PubMed](#) [Entrez Gene](#) [Uniprot](#) [OMIM](#) [Orphanet](#) [Mouse](#) [Rat](#)

**calcium channel, voltage-dependent, L type, alpha 1C subunit**

**Entrez Gene Summary**

This gene encodes an alpha-1 subunit of a voltage-dependent calcium channel. Calcium channels mediate the influx of calcium ions into the cell upon membrane polarization. The alpha-1 subunit consists of 24 transmembrane segments and forms the pore through which ions pass into the cell. The calcium channel consists of a complex of alpha-1, alpha-2/delta, beta, and gamma subunits in a 1:1:1:1 ratio. There are multiple isoforms of each of these proteins, either encoded by different genes or the result of alternative splicing of transcripts. The protein encoded by this gene binds to and is inhibited by dihydropyridine. Alternative splicing results in many transcript variants encoding different proteins. Some of the predicted proteins may not produce functional ion channel subunits. [provided by RefSeq, Oct 2012] [read more](#)

**Orphanet rare diseases**

[Timothy syndrome](#), [Brugada syndrome](#)

**Mouse Genome Database**

Mice homozygous for mutations that inactivate the gene do not survive to term. Selective ablation in beta cells resulted in impaired insulin secretion and systemic glucose intolerance. Heterozygotes were hypoxicative, showed increased anxiety, and poor motor coordination. [read more](#)

[Excel Export](#)

1 < 4 Page 1 of 15 > 200

Group A

Allele frequency

Variations per gene

Genotype  0/0  0/1  1/1  1/2

ROH spans variation  yes/no

Allele frequency

Variations per gene

Genotype  0/0  0/1  1/1  1/2

ROH spans variation  yes/no

Group B

Allele frequency

Variations per gene

Genotype  0/0  0/1  1/1  1/2

ROH spans variation  yes/no

Allele frequency

Variations per gene

Genotype  0/0  0/1  1/1  1/2

ROH spans variation  yes/no

Global Settings

Max Var Frequency

Min Coverage

Min Quality

MaxTargetDist

Transcript Biotype  Protein coding  Other

CDS  UTR  INTRON  SPANNING GENES

PROMOTER  INTERGENIC  COMPOSED

Show Transcripts

CCDS/Refseq  Other

SNP  INSERT  CNV  DELETION

INDEL  ROH

Consequences

Protein structure affected

Strong 5'SS/3'SS effect  Cryptic 5'SS/3'SS activation

Region

Medium 5'SS/3'SS effect

ESS/ESE change  Other or no consequences

dbSNP/dbVar

Genes

Exclude Genes

Cologne Center for  
Genomics

## [Browse Reads](#)

Dataset PRID265:AID1201:ROL\_0091 Region 16:3119298-3119303 Transcript ENST00000528163 Strand 1 Show reads

GAGTTTCGCTGCTCTGTCAAGCTTCTATGTCCTCTTCCAGTCTCACGGAGCCCCAC-GGGGGGACAAGGAGGAGCTGACCCCCAGAATGCTCTGAACCCCATACTCAAATAG	<b>GATA</b>	TGAC	<b>CCAC</b>	A										
0	3119240	3119250	3119260	3119270	3119280	3119290	3119300	3119310	3119320	3119330	3119340	3119350	3119360	3

## Exon Coverage

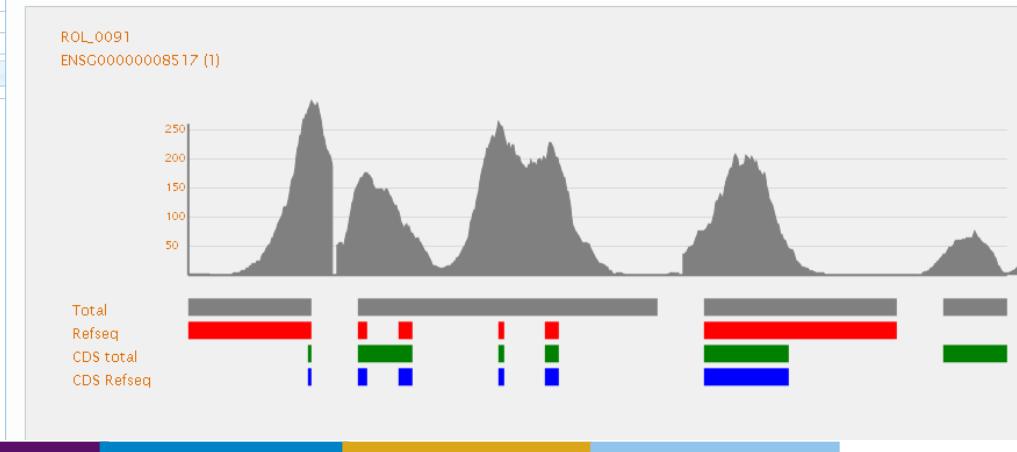
Dataset PRID265:AID1201:ROL\_0091 Gene Ensgene ENSG0000000085

Exon Coverage											
Sample	Ensembl ID	Gene	Transcript	Exon	Chr	Start	Stop	Min_Cov	Ave_Cov	Max_Cov	Perc_Cov
ROL_0091	ENSG000000008517	IL32	ENST00000008180	1	16	3115632	3115688	42	63.51	94	100
ROL_0091	ENSG000000008517	IL32	ENST00000008180	2	16	3115785	3115827	240	274.86	297	100
ROL_0091	ENSG000000008517	IL32	ENST00000008180	3	16	3117378	3117416	0	161.49	176	97.44
ROL_0091	ENSG000000008517	IL32	ENST00000008180	4	16	3117985	3118011	240	254.3	263	100
ROL_0091	ENSG000000008517	IL32	ENST00000008180	5	16	3118181	3118240	194	211.83	230	100
ROL_0091	ENSG000000008517	IL32	ENST00000008180	6	16	3118991	3119820	0	66.06	208	88.92
ROL_0091	ENSG000000008517	IL32	ENST00000325568	1	16	3115298	3115495	0	1.27	3	73.23
ROL_0091	ENSG000000008517	IL32	ENST00000325568	2	16	3115785	3115827	240	274.86	297	100
ROL_0091	ENSG000000008517	IL32	ENST00000325568	3	16	3117378	3117416	0	161.49	176	97.44
ROL_0091	ENSG000000008517	IL32	ENST00000325568	4	16	3117555	3117614	72	86	111	100
ROL_0091	ENSG000000008517	IL32	ENST00000325568	5	16	3117985	3118011	240	254.3	263	100
ROL_0091	ENSG000000008517	IL32	ENST00000325568	6	16	3118181	3118240	194	211.83	230	100

[Home](#)   [Dataset Selection/Filter Options](#)   [Browse Exomes](#)   [Browse Reads](#)   [Gene Coverage](#)   [Download](#)   [Exon Coverage](#)   [Documentation](#)

## Gene coverage

Dataset PRID265:AID1201:ROL\_0091 Gene ENSG00000008517 X-scaling 5 Y-scaling 250 Show coverage



Username   
Password

Login Guest

# varbank

Cologne  
Center for  
Genomics



# **Conclusion**

**WES is extremely successful for disease variant identification in Mendelian diseases → clinical geneticists and pediatricians ask for diagnostic exome-seq!**

- few exomes of affected individuals in consanguineous pedigrees
- 3 or more individuals (affected and unaffected) in dominant families (combination with linkage)
- parents/offspring trios for *de-novo* mutations

**CCG's contribution:**

**>1.500 exomes so far,  
(1/3 monogenic, 1/3 cancer, 1/3 complex diseases)**

# Acknowledgments/Responsibilities

## Lab & Planing

- **Janine Altmüller**  
(project managment)
- **Christian Becker**  
(library prep, sequencing)
- **Elisabeth Kirst**  
(library prep, sequencing)
- **Marek Franitza**  
(library prep, sequencing)

## Computing Center

- **Viktor Achter**  
(HPC administration)
- **Ulrich Lang**  
(resource management)

## Bioinformatics

- **Susanne Motameny**  
(HPC exome pipeline, CNV)
- **Peter Frommolt**  
(HPC, NGSRich, Dindel,  
Tumor/Transcriptome analysis)
- **Kamel Jabbari**  
(SV, CoGIE project data analysis)
- **Amit Kawala**  
(BWA Alignment)
- **Holger Thiele**  
(pipeline planning, data integration,  
functional annotation, trios, database, web  
front end)
- **Wilfried Gunia**  
(server administration, deNovo@HPC)



Username   
Password

Login Guest

# varbank

Cologne  
Center for  
Genomics

