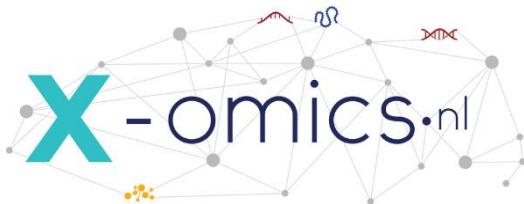




university of
groningen

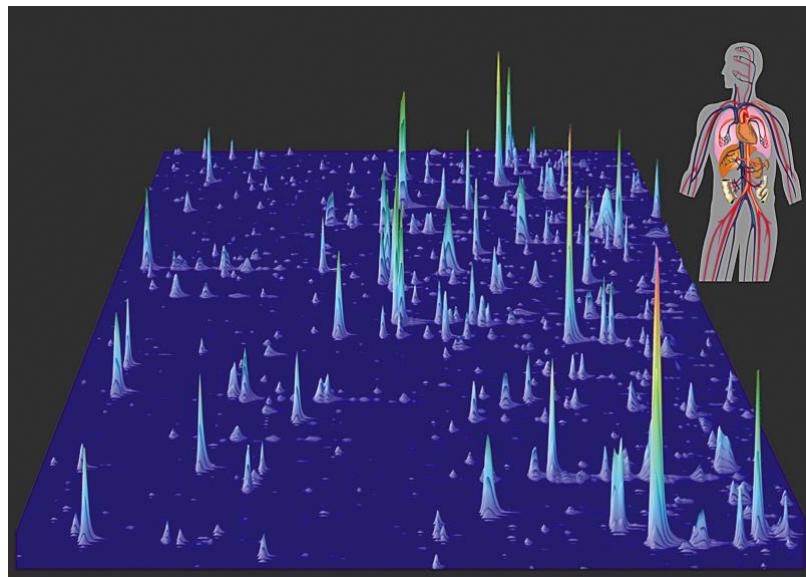


Department of Pharmacy
Analytical Biochemistry

The role of proteogenomics in understanding molecular mechanisms of COPD

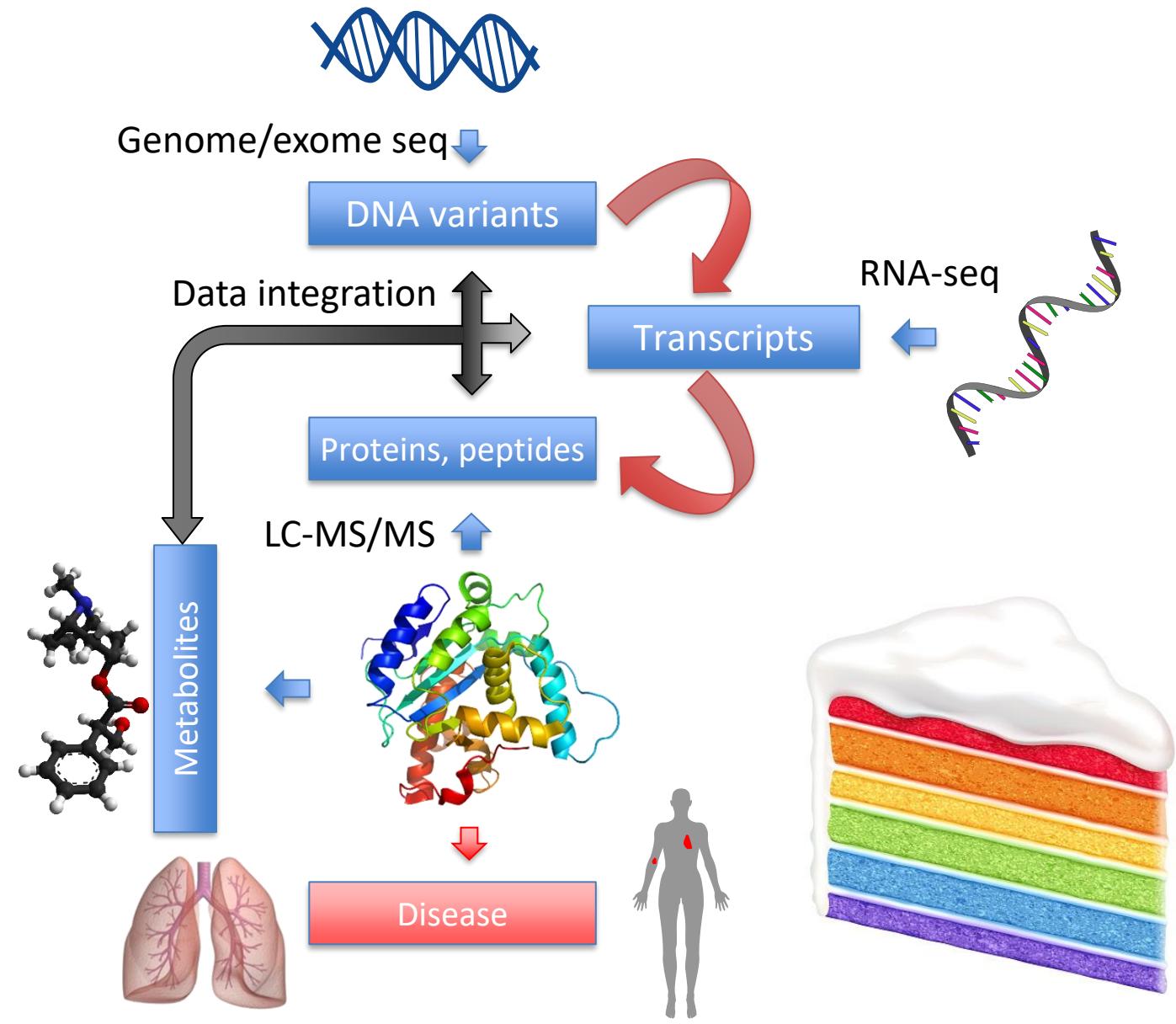
Yanick Hagemeijer, Rainer Bischoff, Victor Guryev, Peter Horvatovich

p.l.horvatovich@rug.nl



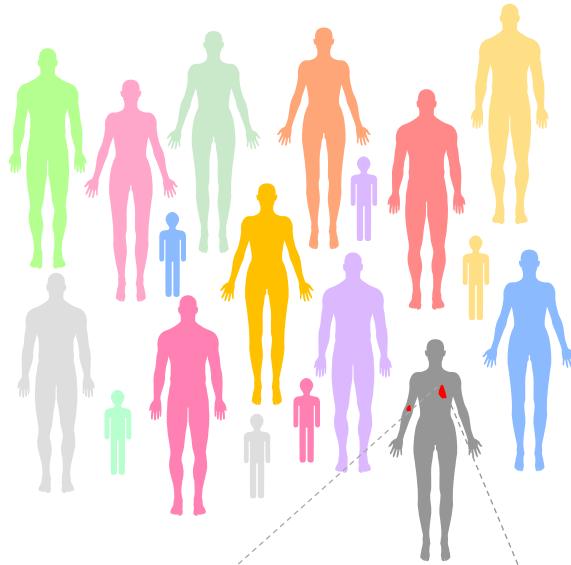
X-Omics Festival, April 12, 2021



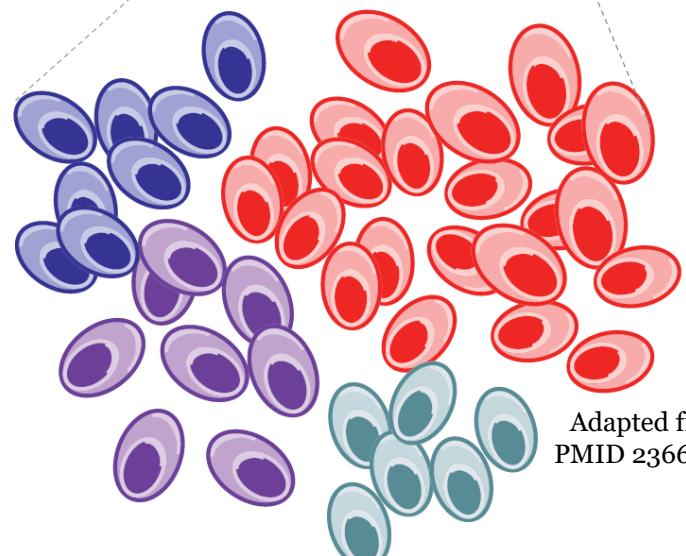


Omics layers:

Phenome
Microbiome
Metabolome
Lipidome
Glycome
Glycoproteome
Kinome
Phosphoproteome
Proteome
Translatome
Transcriptome
Epi-genome
Genome
Metabolome
Lipidome
Acetylome
...



Population genetic variability

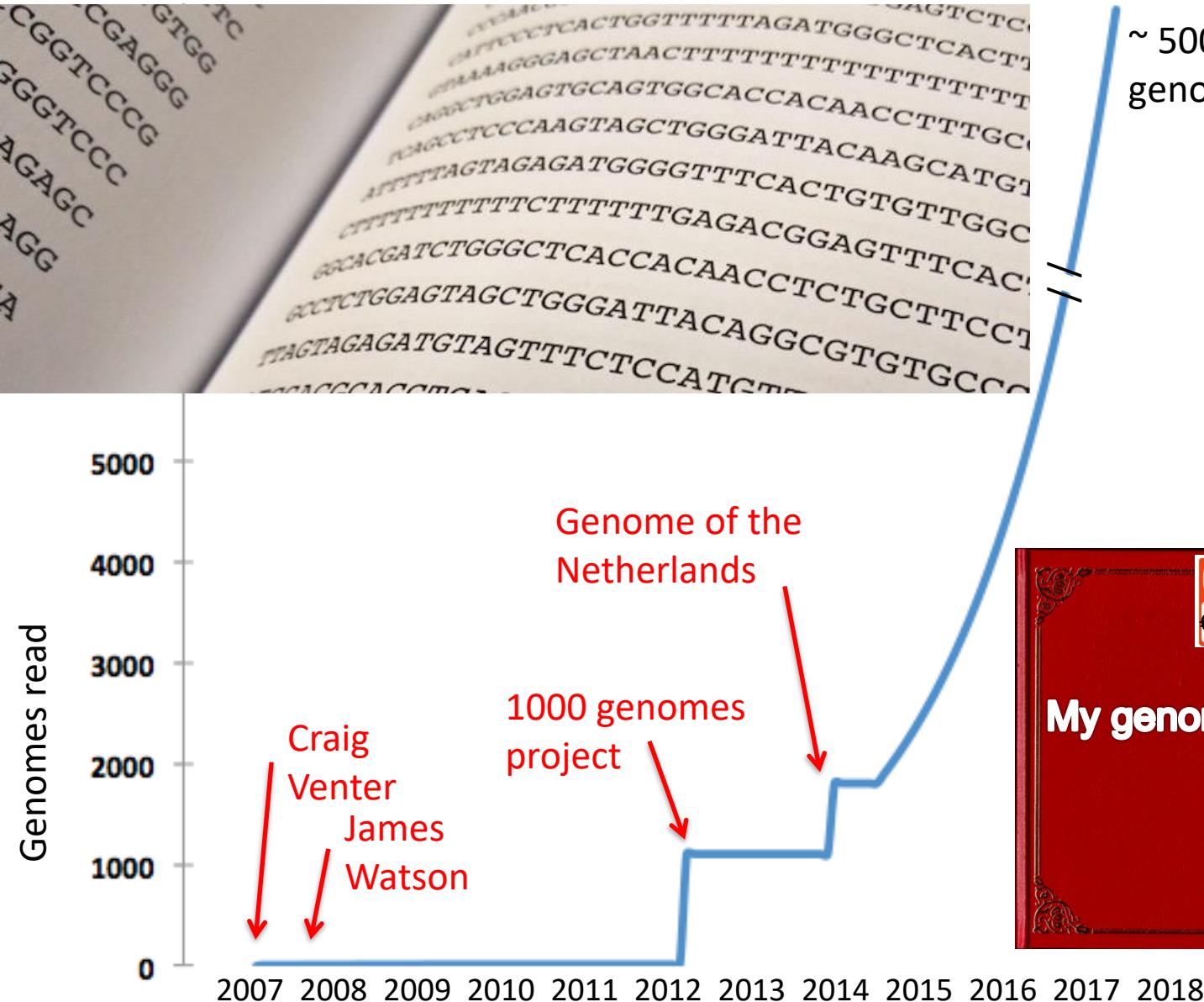


Adapted from
PMID 23664091

Tumor heterogeneity



Era of personal genomics



~ 500 000
genomes

nature
genetics

LETTERS

<https://doi.org/10.1038/s41588-018-0273-y>

OPEN

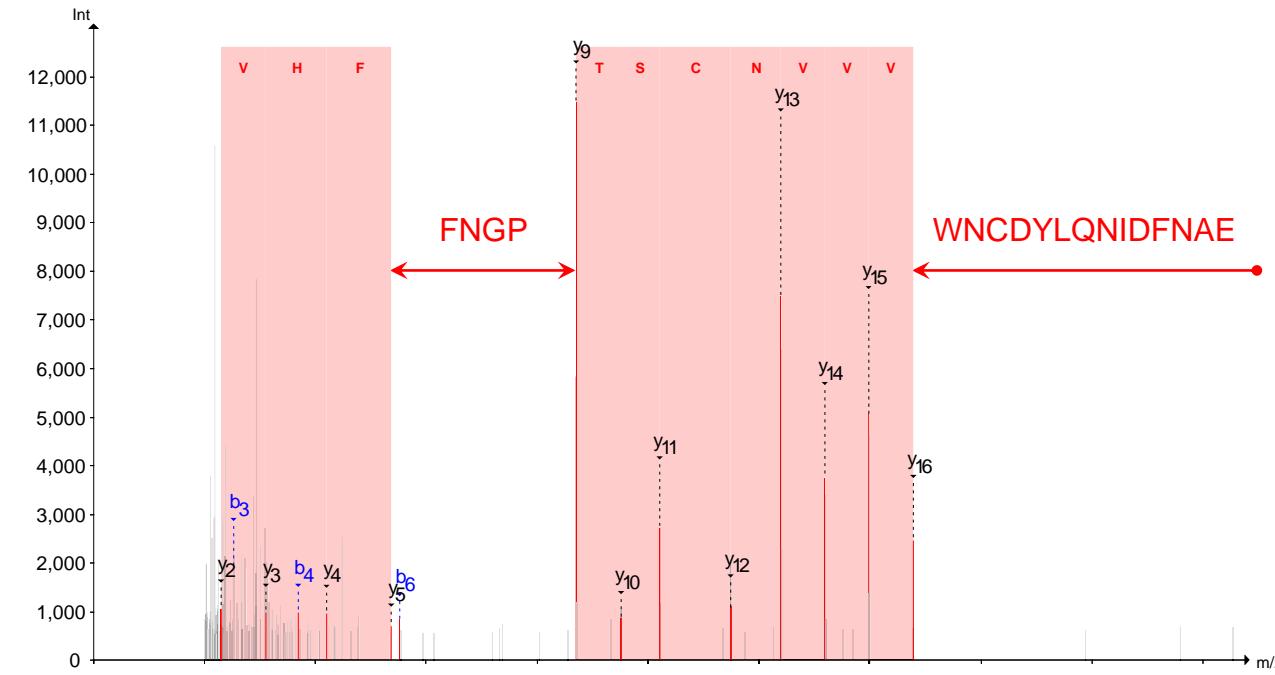
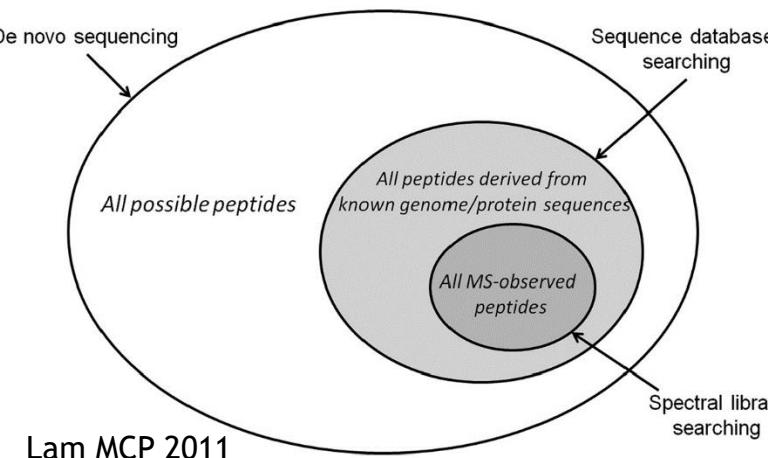
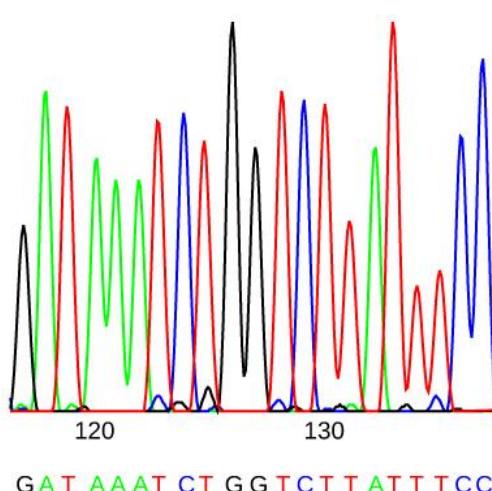
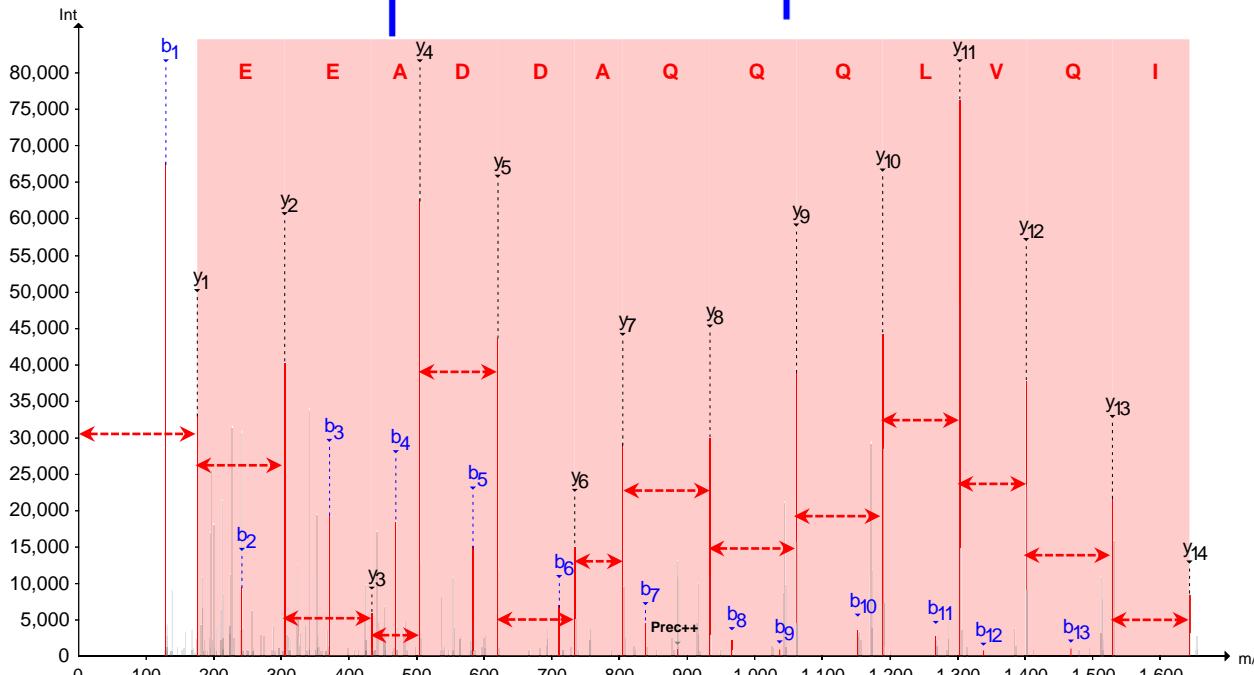
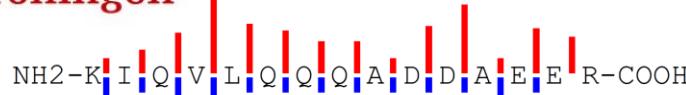
Assembly of a pan-genome from deep sequencing of 910 humans of African descent

We used a deeply sequenced dataset of 910 individuals, all of African descent, to construct a set of DNA sequences that is present in these individuals but missing from the reference human genome. We aligned 1.19 trillion reads from the 910 individuals to the reference genome (GRCh38), collected all reads that failed to align, and assembled these reads into contiguous sequences (contigs). We then compared all contigs to one another to identify a set of unique sequences representing regions of the African pan-genome missing from the reference genome. Our analysis revealed 296,485,284 bp in 125,715 distinct contigs present in the populations of African descent, demonstrating that the African pan-genome contains -10% more DNA than the current human reference genome. Although the functional significance of nearly all of this sequence is unknown, 387 of the novel contigs fall within 315 distinct protein-coding genes, and the rest appear to be intergenic.

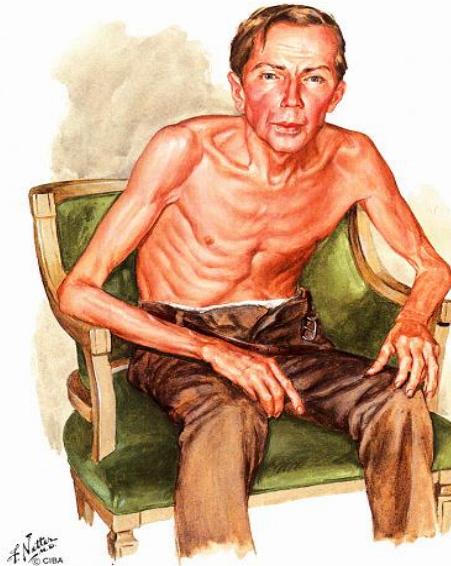
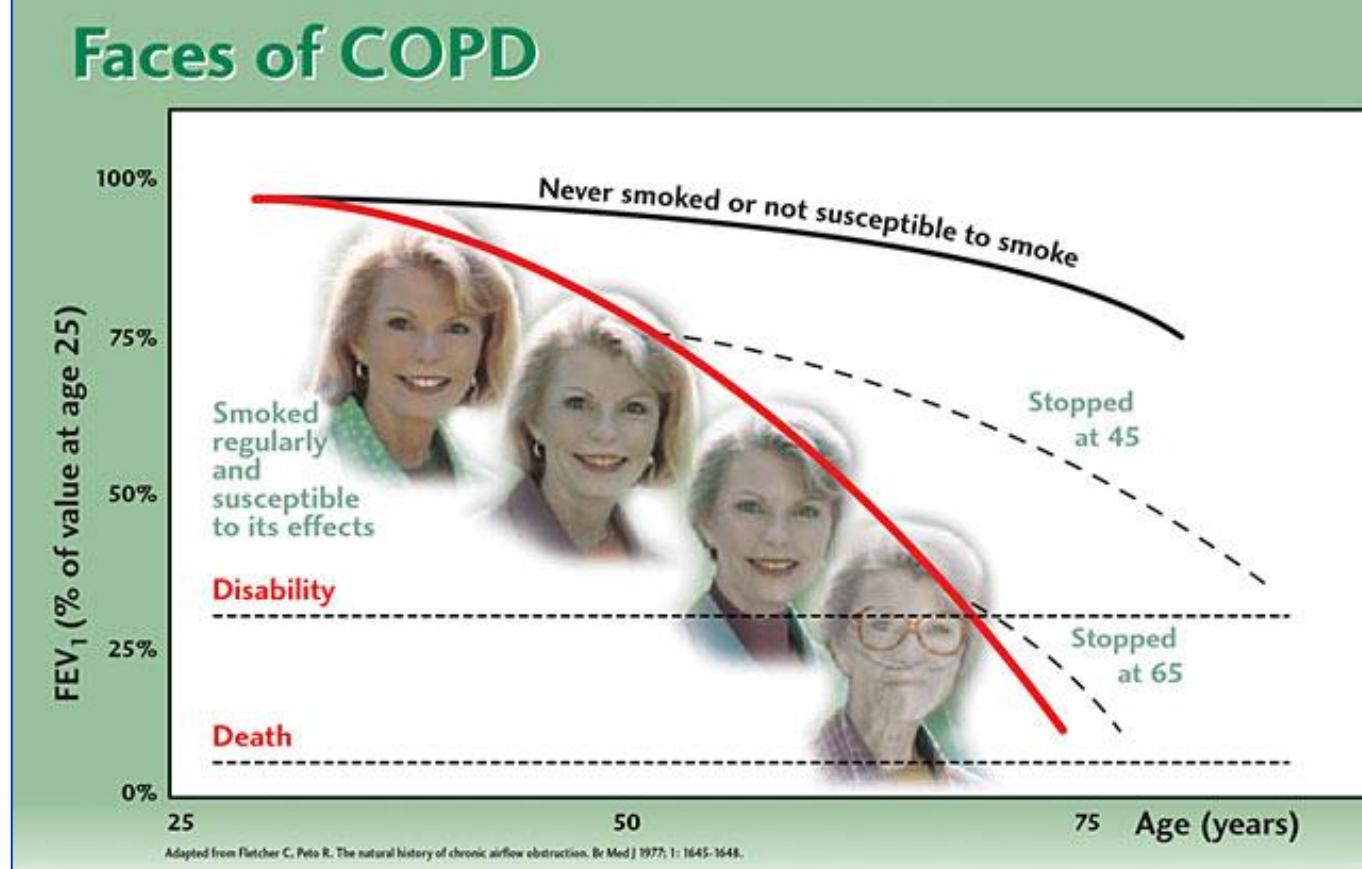


university of
groningen

Difficulty in identifying variants *de novo* from MS/MS spectra



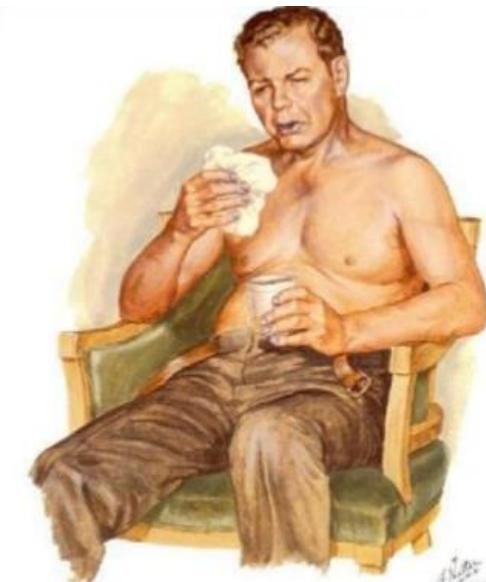
- Best identification method: database search (DBS)
- DBS requires list of protein sequences expected to be present in the sample
- Canonical sequences Swissprot and Uniprot (Ensembl) is used in common proteomics workflows
- In public databases low number of variants (20, 80 and 30 k proteins) are present.



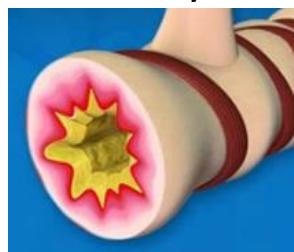
- 20% of the smokers develop COPD, more than 200 millions persons have COPD
- Progressive loss of lung function with a large impact on the quality of life
- Insufficient insight in the molecular mechanisms of COPD
- Limited therapeutic options



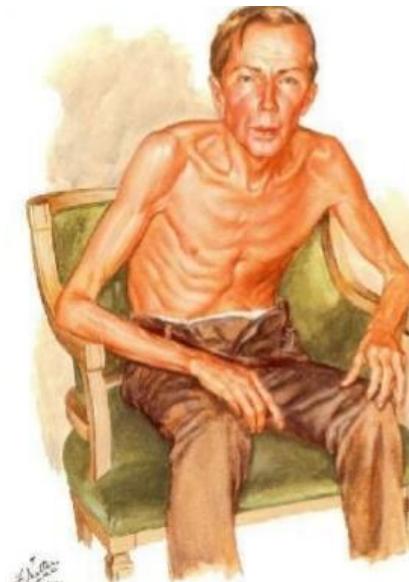
chronic bronchitis



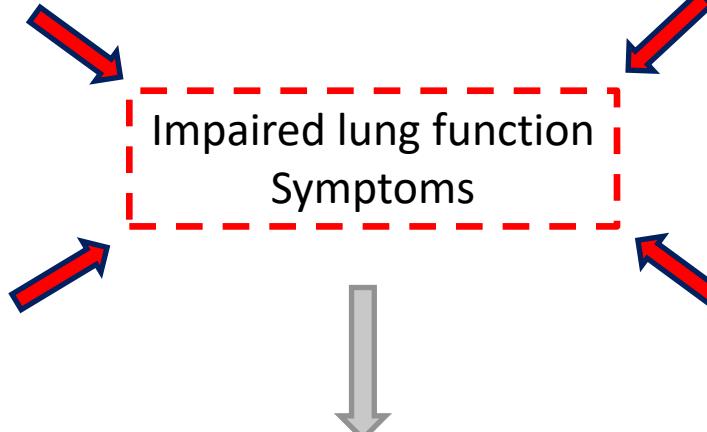
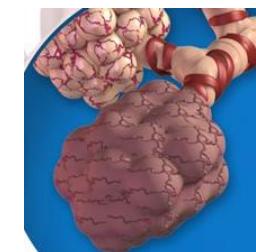
airways



emphysema



alveoli



Today treatment “**One size fits all**”



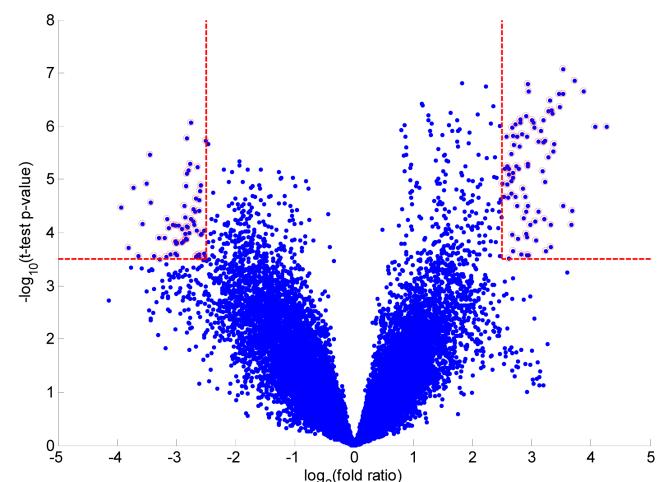
Study design

20/18 lung tissues

10 COPD stage IV (8 female/2 male)
10/8 ex-smokers control (6/4 female/4 male)

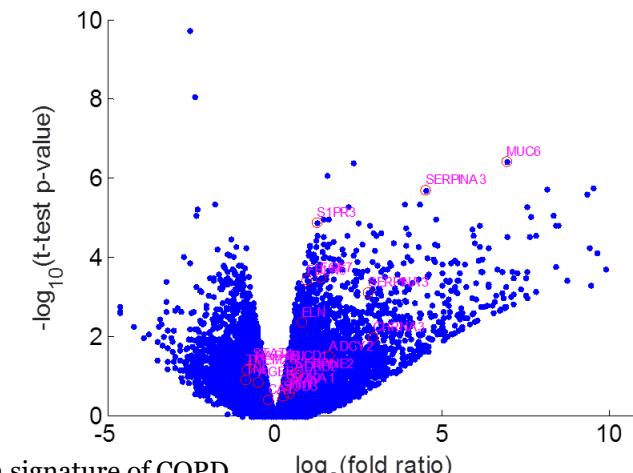
proteins

Orbitrap Q-Executive+ 1D-LC-MS/MS 2D-LC-MS/MS



mRNA sequences

Illumina
20 million sequencing depth
Polyadenylated mRNA fraction



COPD

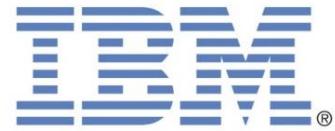
Control



university of
groningen



Proteogenomics data integration workflow



mRNA-Seq
mRNA-Seq pre-processing

Proteomics
MS1 quantification

Proteomics
peptide/protein identification

mRNA raw sequence
Illumina

FastQC
quality control

Trimmomatic
sequence trimming

- Annotated proteins
- SAAVs
- Splice isoforms
- indels
- New gene models
- New transcripts

STAR
alignment and annotation

TopHat2
alignment and annotation

Cufflinks
mRNA assembly and quantification

StringTie
mRNA assembly and quantification

TransDecoder
Translated protein sequence

mRNA quantitative table
Fasta format/FPKM

SearchGUI/PEAKS
peptide/protein identification

PeptideShaker/PEAKS
validation of peptide/protein identification

LC²-MS/MS data
Orbitrap QE+

Grid
smoothing, resampling

Centroid
peak detection and quantification

Warp2D
time alignment

MetaMatch
peak matching

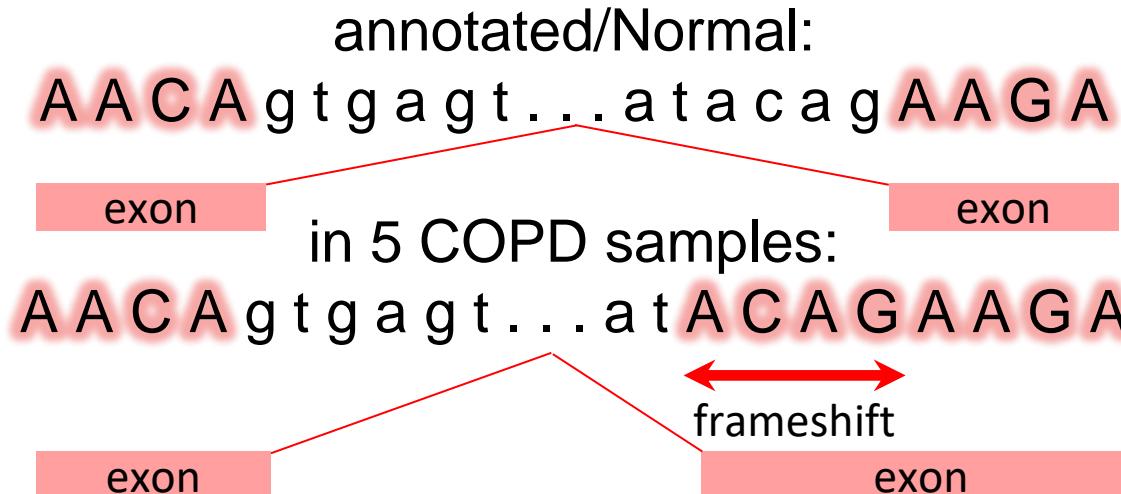
nextflow nf-core



Input files
mRNA-Seq and LC-MS/MS

Output files
Protein/mRNA and quantities

Suits F, Hoekman B, Rosenling T, Bischoff R, Horvatovich P., Threshold Avoiding Proteomics Pipeline,
Anal Chem., 2011, **83**(20):7786-94.



One reason why transcript is present, and protein is absent

translated protein sequence of oncostatin M receptor gene without mutation

5'3' Frame 1

Met ALFAVFQTTFLTLLSLRTYQSEVLAERLPLTPVSLKVSTNSTRQLHLQWTVHNLPYHQELK Met VFQIQISRIETSNVIWVGNYSTTVKWNQVLHWSWESELPLECATHFVRIKSLVDDAKFPEPNFWSNWSSWEEVSVQDSTGQDILFVFPKDKLVEEGTNVTICYVSRNIQNNVSCYLEGKQIHGEQLDPHTAFNLNSVPFIRNKGNTIYCEASQGNVSEG Met KGIVLFVSKVLEEKDFSCETEDFKTLHCTWDPGTDALGWSKQPSQSYTLFESFSGEKKLCTHKNWCNWQITQDSQETYNTFTLIAENYLRKRSVNILFNLTHERVYL Met NPFSVNFENVNATNAI Met TWKVHSIRNNFTYLCQIELHGEVK Met Met QYNVSIKVNGEYFLSELEPATEY Met ARVRCADASHFWKSEWSGQNF TTL EAAPSEAPDVWRIVSLEPGNHTVTLFWKPLSKLHANGKILFYNVVVENLDKPSSSELHSIPAPANSTKLILDRCSYQICVIANNVGA SPASVIVISADPENKEVEEERIAGTEGGFSLSWKPQPGDVIGYVVWDWCHTQDVLGDFQWKNVGPNTSTVISTDAFRPGVRYDFRINYGLSTKRIACLLEKKTGYSQELAPSNDPHVLVDLTSHSFTLSWKDYSTESQPGFIQGYHVYLVSKARQCHPRFEKAVALSDGSECCKYKIDNPEEKALIVDNLKPESFYEFFITPFTSAGEGPSATFTKVTTPDEHSS Met LIHILLP Met VFCVLLI Met V Met CYLKSQWIKETCYPDI PDPYKSSILSLIKFKENPHLII Met NVSDCIPDAIEVVSKEGTQIQLGTRKSLTETELTKPNYLYLLPTEKHNHSGPGPCICFENLTYNQ AASDSGSCGHVPVSPKAPS Met LGL Met TSPENVLKALEKNY Met NSLGEIPAGETSLNYVSQLASP Met FGDKDSLPTNPVEAPHCSEY K Met Q Met AVSLRLALPPPTENSSSITLLDPGEHYC Stop

with mutation

5'3' Frame 1

Met ALFAVFQTTFLTLLSLRTYQSEVLAERLPLTPVSLKVSTNSTRQLHLQWTVHNLPYHQELK Met VFQIQISRIETSNVIWVGNYSTTVKWNQVLHWSWESELPLECATHFVRIKSLVDDAKFPEPNFWSNWSSWEEVSVQDSTGQDILFVFPKDKLVEEGTNVTICYVSRNIQNNVSCYLEGKQIHGEQLDPHTAFNLNSVPFIRNKGNTIYCEASQGNVSEG Met KGIVLFVSKVLEEKDFSCETEDFKTLHCTWDPGTDALGWSKQPSQSYTLFESFSGEKKLCTHKNWCNWQITQDSQETYNTFTLIAENYLRKRSVNILFNLTHERVYL Met NPFSVNFENVNATNAI Met TWKVHSIRNNFTYLCQIELHGEVK Met Met QYNVSIKVNGEYFLSELEPATEY Met ARVRCADASHFWKSEWSGQNF TTL EAAPSEAPDVWRIVSLEPGNHTVTLFWKPLSKLHANGKILFYNVVVENLDKPSSSELHSIPAPANSTKLILDRCSYQICVIANNVGA SPASVIVISADPENNERG Stop GRKNCRHRGWILSVLETPTWRCYRLCCGLV Stop PYPGCAR Stop FPVEECRSQYHKHSH Stop HRCF Stop ARSSI Stop LQNLWV1YKKDCFLIREKNRILSGTCSRQPSRAGGYIDIPLLHSELERLLY Stop ISTWFYTRVPCCLSEI1QGEAVPPTI Stop KGSSFRWFR Met LQIQN Stop QPGRKGIDCGQPKARILL Stop VFHHSIH Stop CW Stop RPQCYVHEGHDSG Stop TLLDADSYPTAHGF LRLAHHGHVLLKEKSVQDQGDLLS Stop HP Stop PLQEQQHPVINKIQGEPSPNNECQ Stop LYPRCY Stop SCKQARRDKDTVPRH Stop EVTH RNRVD Stop A Stop LPLSPSNRKESLWPWPLHLF Stop ELDL Stop PGSF Stop LWLLWPCSSIPKSPKYAGTNNDLT Stop KCTKGTRKKLHELP GRNPSWRNKFEELCVPGFTHVWRQGQSPNPKPSRGTTLFRV Stop NANGSLPASCLASPDRE Stop QPLLNYPFRSR Stop TLLL

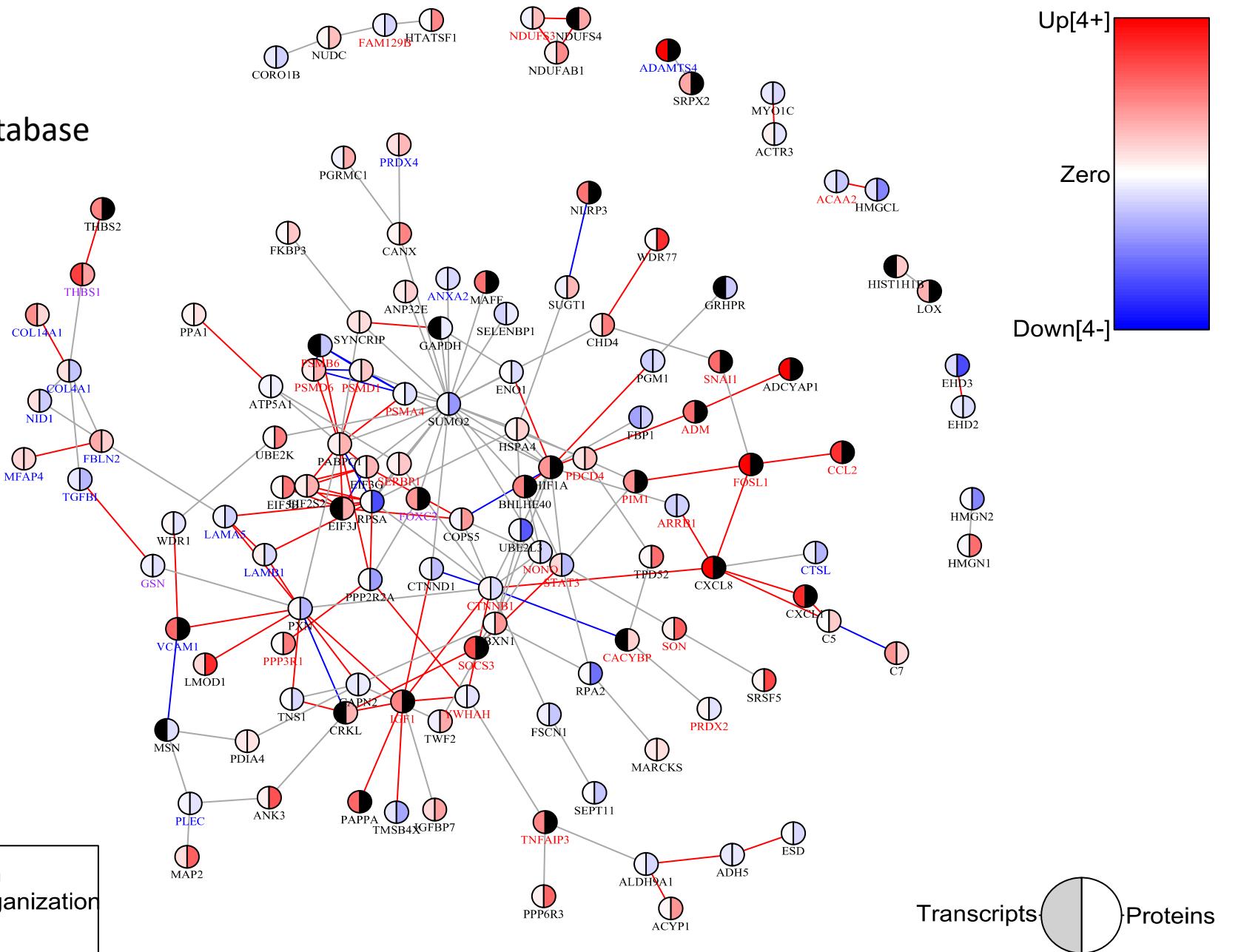
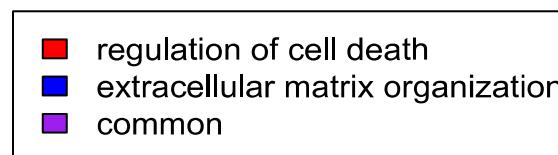


STRING protein-protein interaction network

STRING edges: experiments and database

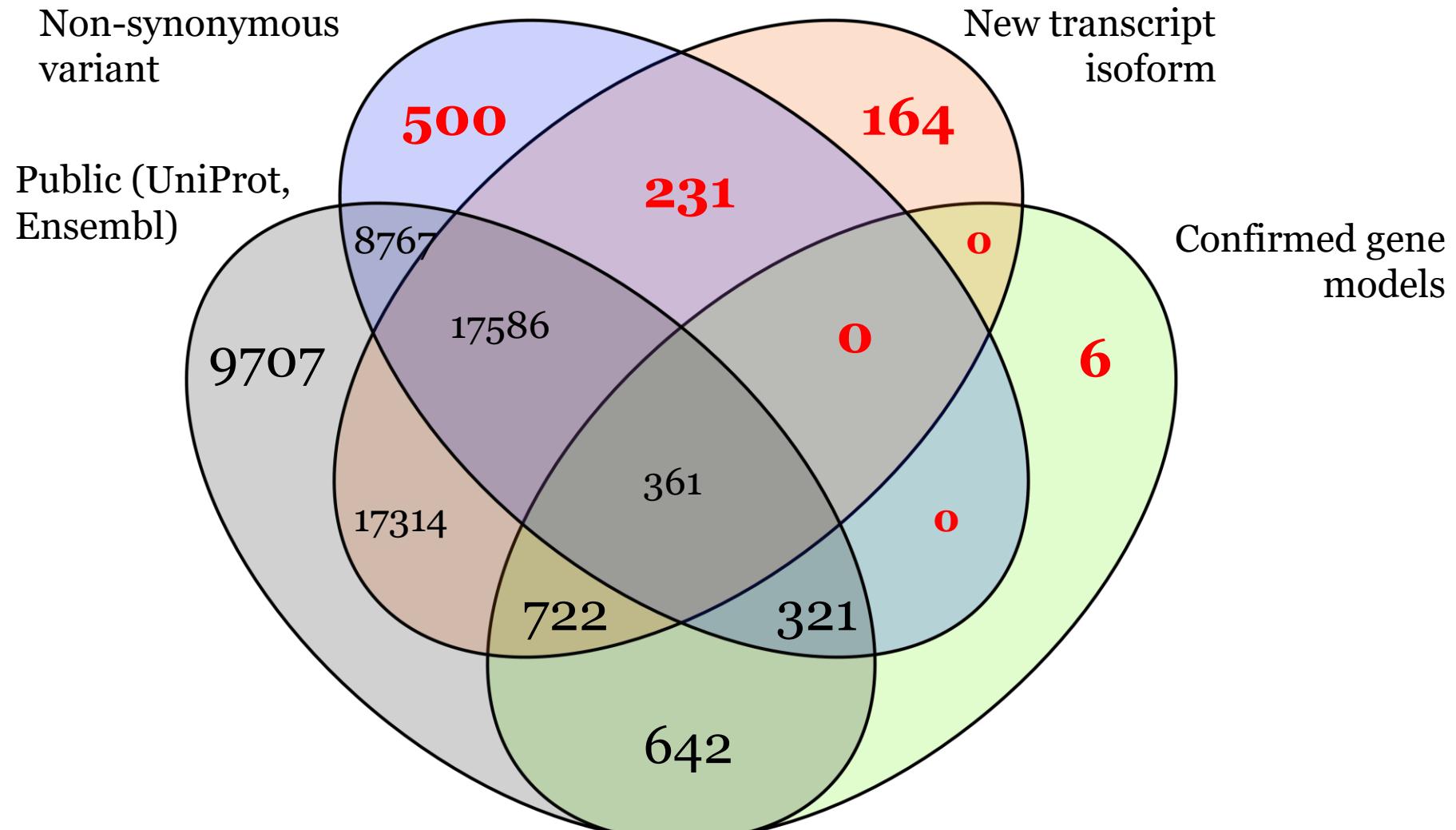
- experiments
— database
— both

Proteomics +
transcriptomics
FDR 0.01
Fixed layout





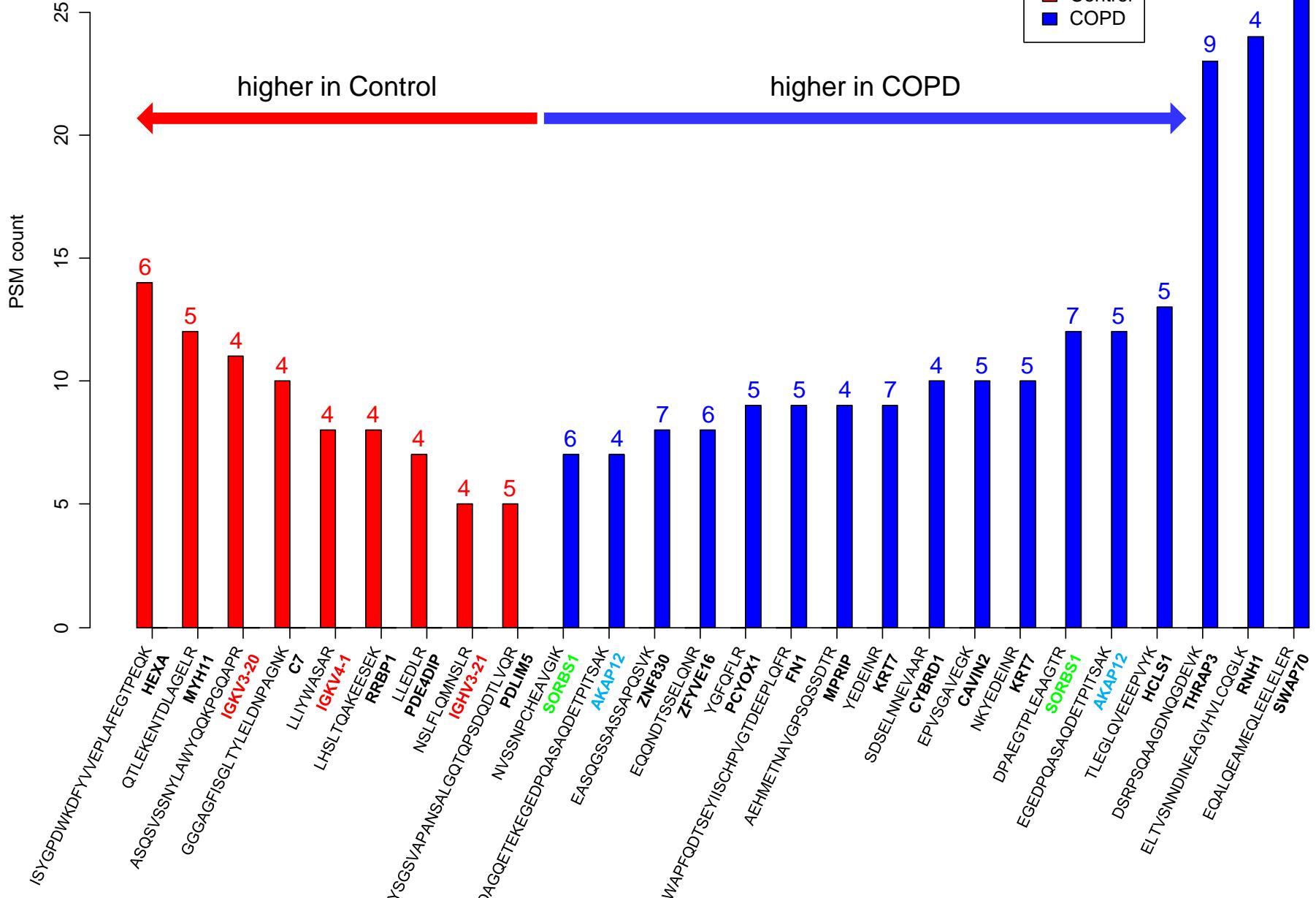
Peptide identification (combined dataset)



901 peptides are not in curated public databases



PSM count of novel peptides present only in COPD or Control

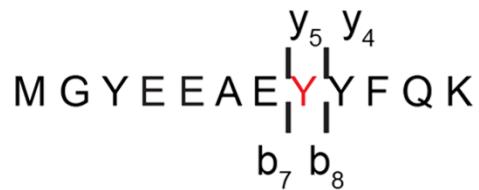




Group only “Novel peptides” identification PEAKS score and SpetcrumAI test

Peptides only present in COPD

Peptide sequence	Quality score	Gene	Effect	Ion support	Flanking ion support
ADSQDAGQETEKEGEDDPQASAQDETPITSAK	134.83	AKAP12	E1600D amino acid substitution	x	x
DSRPSQAAGDNQGDEVK	132.86	THRAP3	A201V amino acid substitution	x	x
TGQEALSQTTISWAPFQDTSEYIISCHPVGTDEEPLQFR	131.44	FN1	Native	NA	NA
AEHMETNAVGPSQSSDTR	129.43	MPRIP	P327Q amino acid substitution	x	x
EGEDDPQASAQDETPITSAK	129.34	AKAP12	E1600D amino acid substitution	x	x
DPAEGTPLEAAGTR	118.25	SORBS1	Splice variant, exon extension	NA	NA
ELTVSNNDINEAGVHVLQGLK	115.80	RNH1	R188H amino acid substitution	x	x
TLEGLQVEEEPVYK	106.37	HCLS1	E361K amino acid substitution	x	x
EASQGSSASSAPQSVK	105.60	ZNF830	H99Q amino acid substitution	x	x
SDSELNNEVAAR	103.81	CYBRD1	S266N amino acid substitution	x	x
EQALQEAMEQLEELELER	100.83	SWAP70	Q505E amino acid substitution	x	x
EQQNDTSSELQNR	91.55	ZFYVE16	I192T amino acid substitution	x	x
YGFQFLR	81.46	PCYOX1	S149F amino acid substitution	x	x
NKYEDEINR	70.58	KRT7	H186R amino acid substitution	x	x
EPVSGAVEGK	62.17	CAVIN2	Q174E amino acid substitution	x	x
NVSSNPACHEAVGIK	56.53	SORBS1	Splice variant, exon extension	NA	NA
YEDEINR	49.80	KRT7	H186R amino acid substitution	x	x

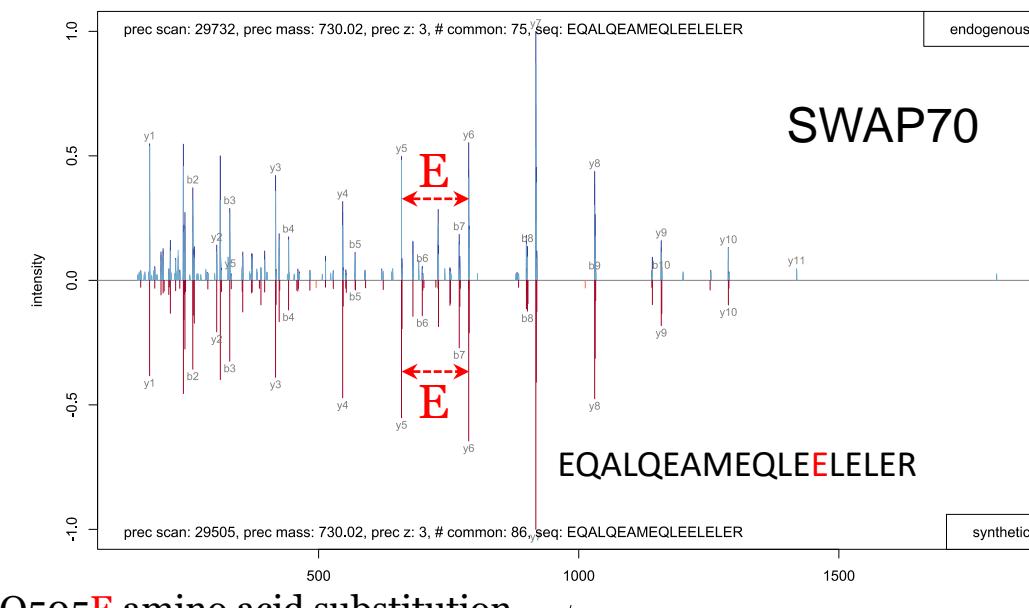
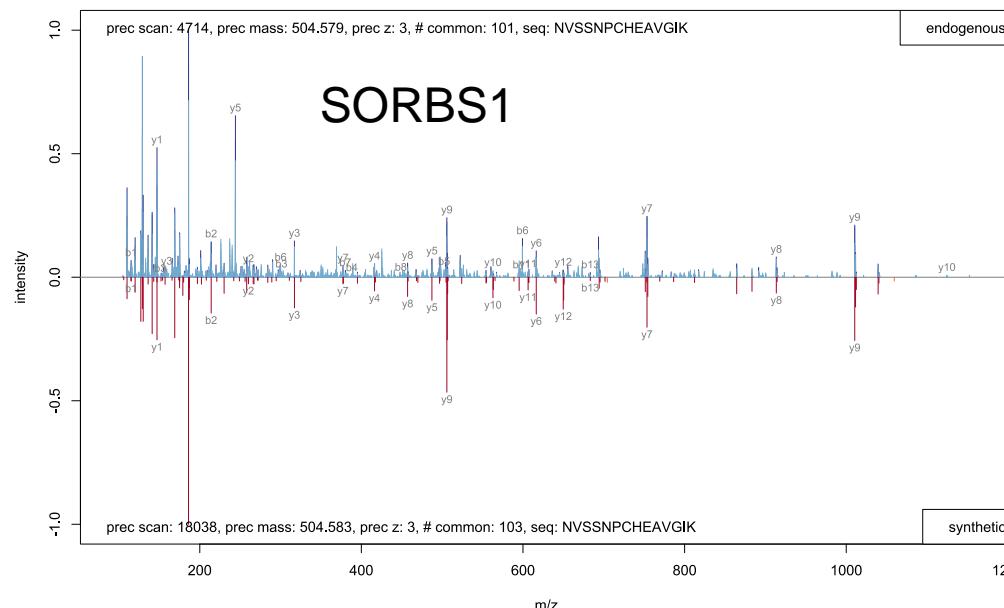
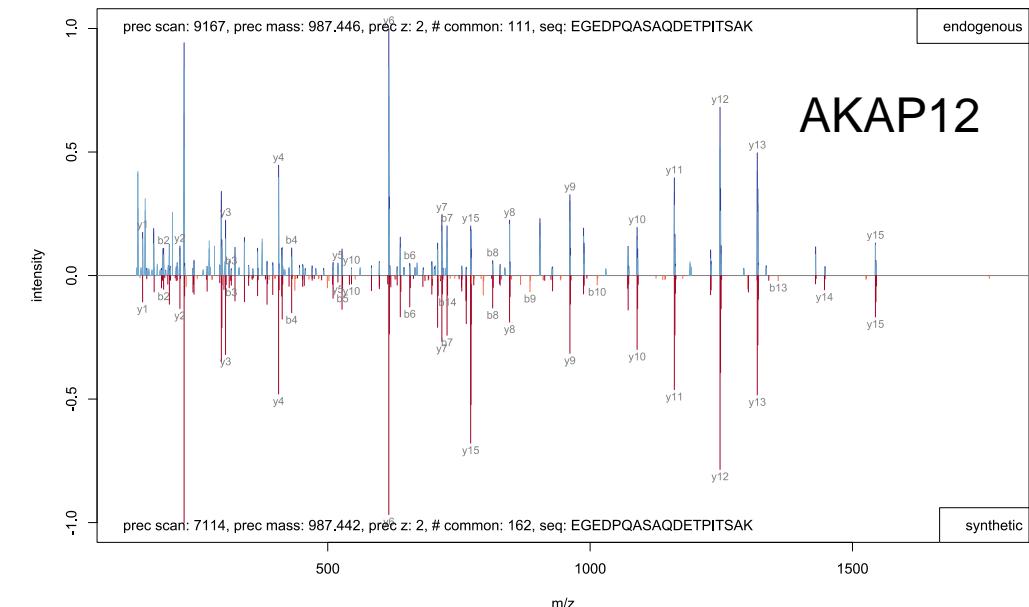
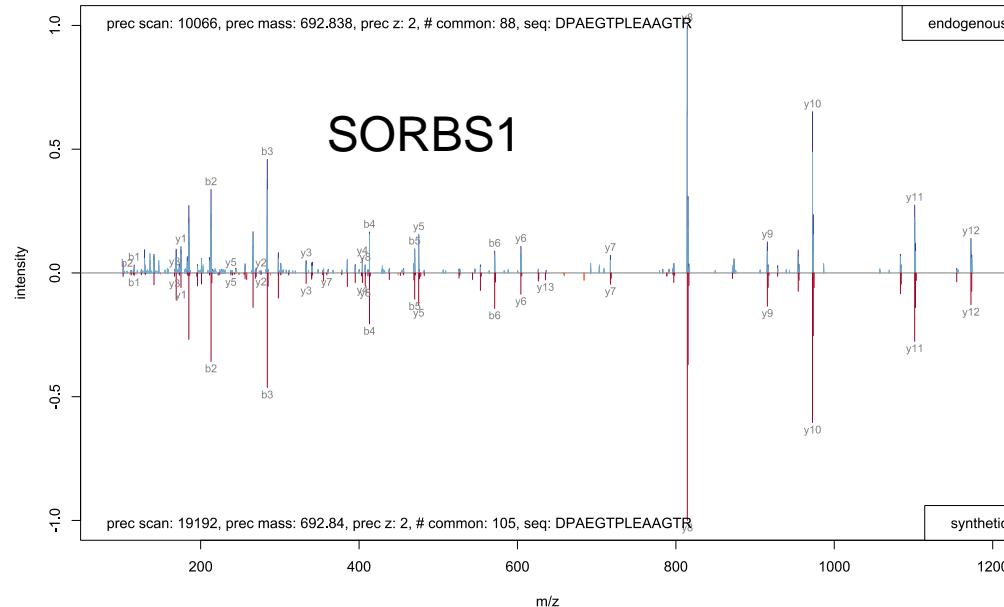


Peptides only present in Control

Peptide sequence	Quality score	Gene	Effect	Ion support	Flanking ion support
ISYGPDWKDFYVVEPLAFEGTPEQK	120.51	HEXA	I447V amino acid substitution	x	x
ISNSAAYSGSVAPANSALGQTQPSDQDTLVQR	110.84	PDLIM5	T410A amino acid substitution	x	x
LHSLTQAKEESEK	110.58	RRBP1	L1043H amino acid substitution	x	-
GGGAGFISGLTYLELDNPAGNK	108.34	C7	S389T amino acid substitution	x	x
ASQSVSSNYLAWYQQKPGQAPR	105.88	IGKV3-20	S52N amino acid substitution	x	x
QTLEKENTDLAGELR	77.52	MYH11	A1241T amino acid substitution	x	x
NSLFLQMNSLR	76.68	IGHV3-43	Y99F amino acid substitution	x	x
LLIYWASAR	68.41	IGKV4-1	T79A amino acid substitution	x	x
LLEDLR	33.62	OTOA/PDE4DIP	Native	NA	NA



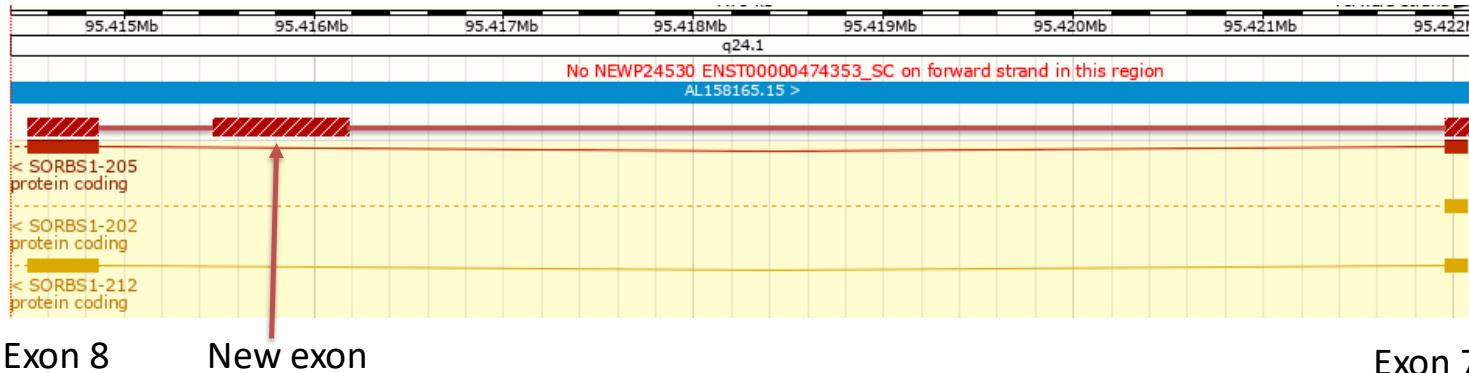
Validation of non-reference peptides with synthetic peptides MS/MS spectra





Two peptides uniquely mapping to a novel exon of SORBS1 gene

SORBS1 is encoded on minus strand of chr10 (band 10q24.1), between 95.31 and 95.56 Mbp (gene length: 249.64 kb)



>NEWP24530

```
MSSECDGGSKAVMNGLAPGSNGQDKATADPLRARSISAVKIIIPVKTVKNASGLVLPNDMD
LTKICTGKGAVTLRASSSYRETPSSSPASPOETRQHESKPGLPEPSSADEWRLSSSADA
NGNAQPSSLAAGYRSVHPNLPDKSQDATSSAAQPEVIVVPLYLVNTDRGQEQTARPP
TPLGPLGCVPТИPATASAASPLTFPTLDDFIPPHLQRWPHHSQPARASGSFAPISQTPPS
FSPPPPLVPPAPEDLRRVSEPDLTGAVSSTVLSPPRPLLQKDRAFWQSPTIHNTYKDSL
YLSSPKPYVPLGTPRQQNPSQPQPISVLLAAGSAPKGVCPGSLLPDSTFPSASSQPQQR
YAATRTVYHKNVSSNPCHAEVGIKKVSSLYVPCLSNNICLAASENSSRVARDPAEGTPLE
AAGTRAPAPGLVSRTAGTGKPPPAPPDPPLKFDIRKDAVNREGESPLSLGTQASFDPDVRP
PVLGPRVTSDPENRKSKESYLLQPSYPAKDSSPLLNEVSSSLIGTDSQAFPSVSKPSSAY
PSTTIVNPTIVLLQHNREQQKRILSSLSDPVSERVGEQD SAPTQEKPSPGKAIEKRAK...
```

Exons 1-7

New exon

Exons 8 and further

Protein domains of SORBS1 splice variants

native SORBS1



novel SORBS1 splice variant



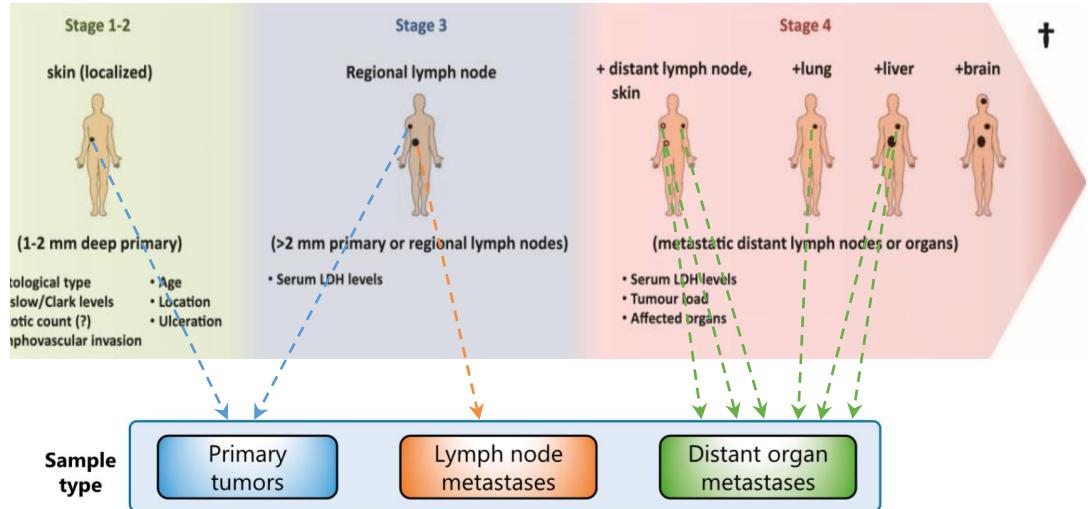


university of
groningen

Cancer moonshot for melanoma

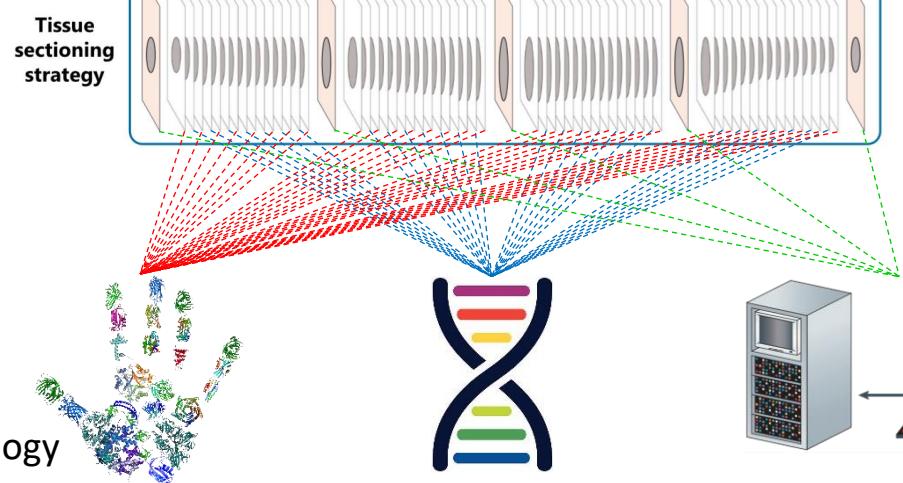


Clinically well
characterized
patients



Biobanks

Tumor sections

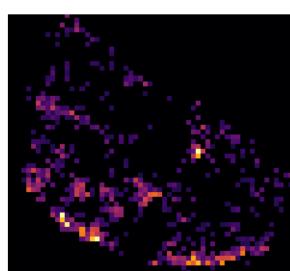
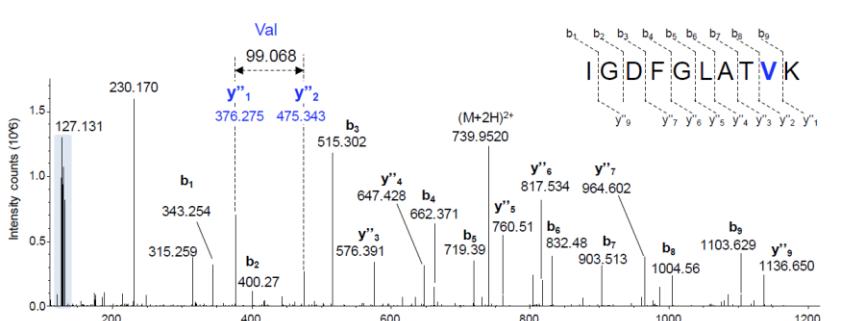
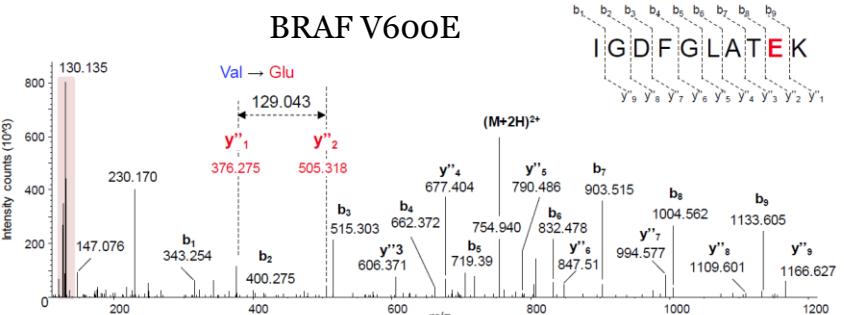
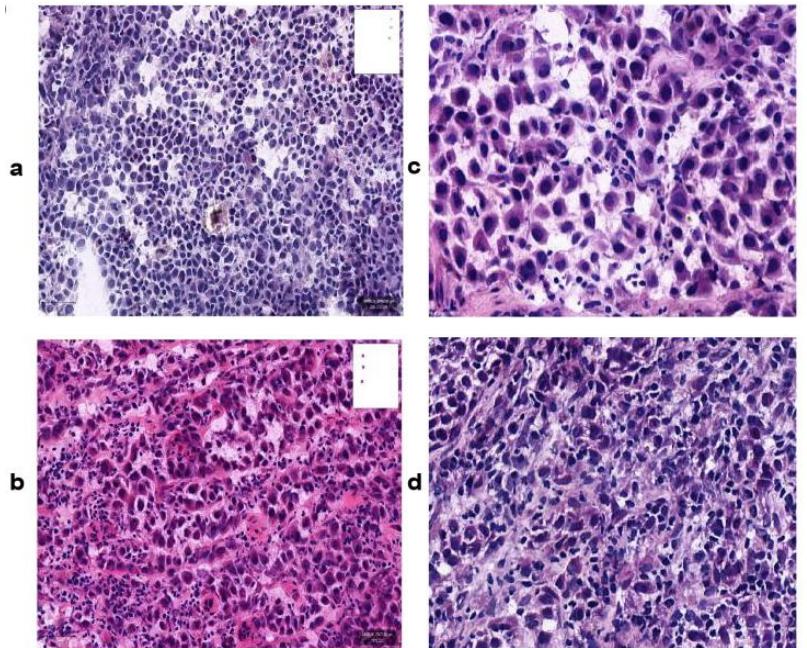


Multi-omics
profiles
& digital pathology

Proteomics

Genomics

Digital pathology





university of
groningen

ERIBA

Yanick Hagemeijer
Ekaterina Ovchinnikova
Victor Guryev



university of
groningen

Alejandro Sanchez Brotons
Karel Gerbrands
Ana Ciconelle
Hjalmar Permentier
Rainer Bischoff



Fundings and Collaborations



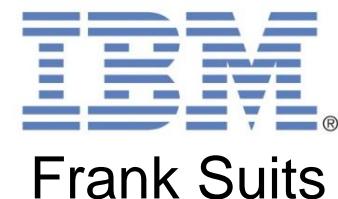
umcg



Corry-Anke Brandsma
Maarten van den Berge
Wim Timens
Karin Wolters
Dirkje Postma
Gyorgy Halmos
Renee Vehoeven



Maria Yakovleva
Melinda Rezeli
Jonatan Eriksson
Thomas Fehniger
György Markó-Varga



Frank Suits



ERIBA