

QC metrics

Somatic and germline variant calling

Outline

- **Quality Control - QC**
- Metrics
- Tools
- Correlations
 - Plots
- Quiz
- **Variant calling**
- Tumor purity and clonality
- Somatic mutation vs. germline mutation vs. germline polymorphism
- Germline variant calling
 - Methodology
 - Tools
 - Quiz
- Somatic variant calling
 - Methodology
 - Tools
 - Quiz
- Variant annotation (VEP)
- File format
- Manual curation
- Microsatellite instability/hypermutation

QC metrics

Quality control metrics

- Was the sequencing successful or not?
- Many steps can go wrong
 - Extraction of DNA, DNA input amount
 - Library prep, e.g. PCR amplification
 - Capture
 - Sequencing
 - Demultiplexing
- Important with quality control metrics
- Is data quality “good enough”?
 - Requirements can vary a lot depending on the experiment

Quality control metrics

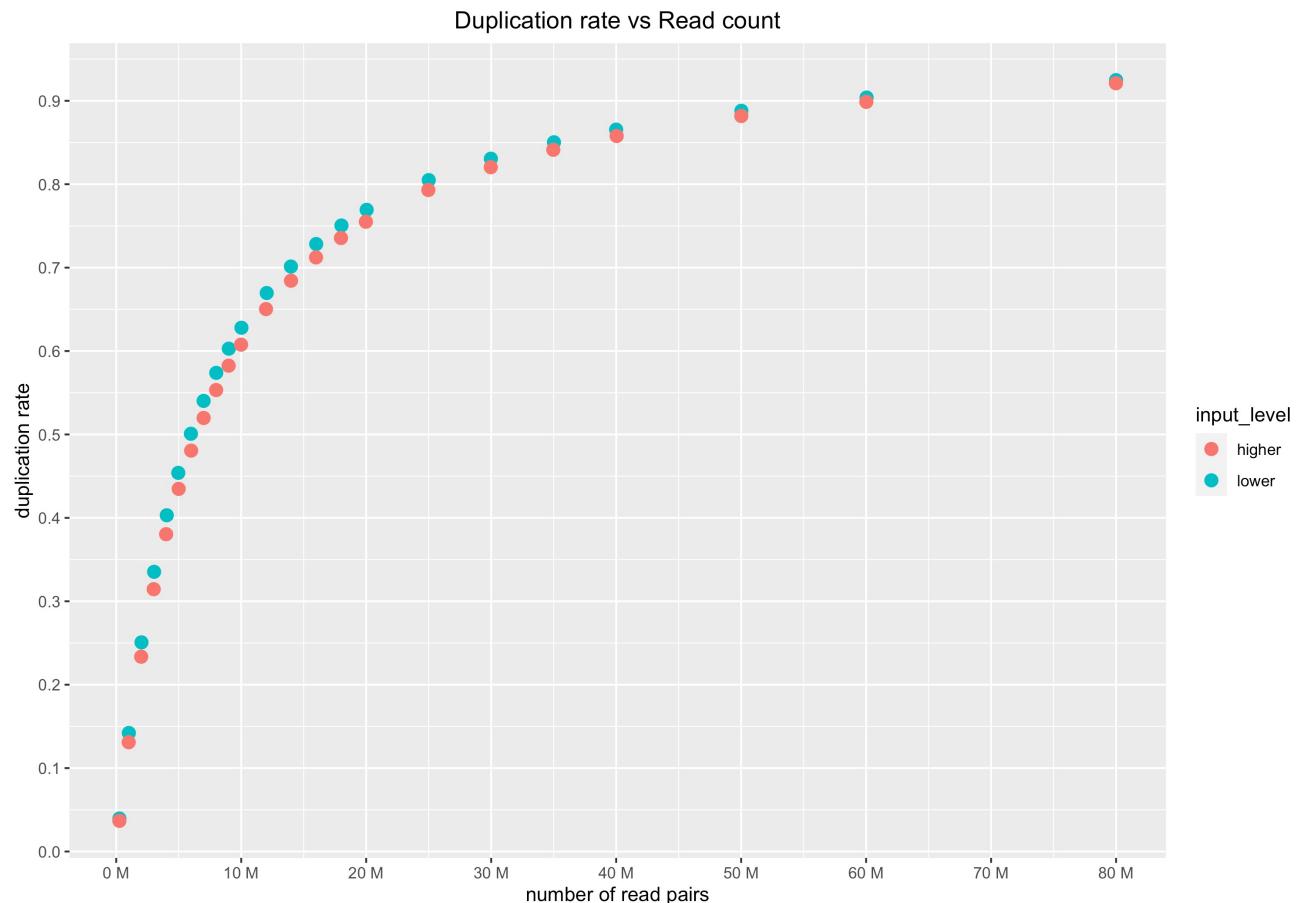
- The most important metrics
 - **Coverage** (after deduplication) – the average number of (unique) reads covering the targeted regions, also known as depth
 - **Read count** – total number of reads for a sample
 - **Duplication rate** – what fraction of all reads where duplicates (not unique)
 - **Fold enrichment** – how much more the targeted regions are amplified compared to non-targeted regions, x-fold
 - **Contamination** – DNA from another source

Quality control metrics

- Example of tools for QC
 - Picard CollectHsMetrics
 - Coverage, fold enrichment
 - Picard MarkDuplicates
 - Read count, duplication rate
 - GATK (v ≤3) ContEst
 - Contamination
- Input: bam files
- Output: txt tables,
 - can be parsed in e.g. R for plotting, summary tables etc.

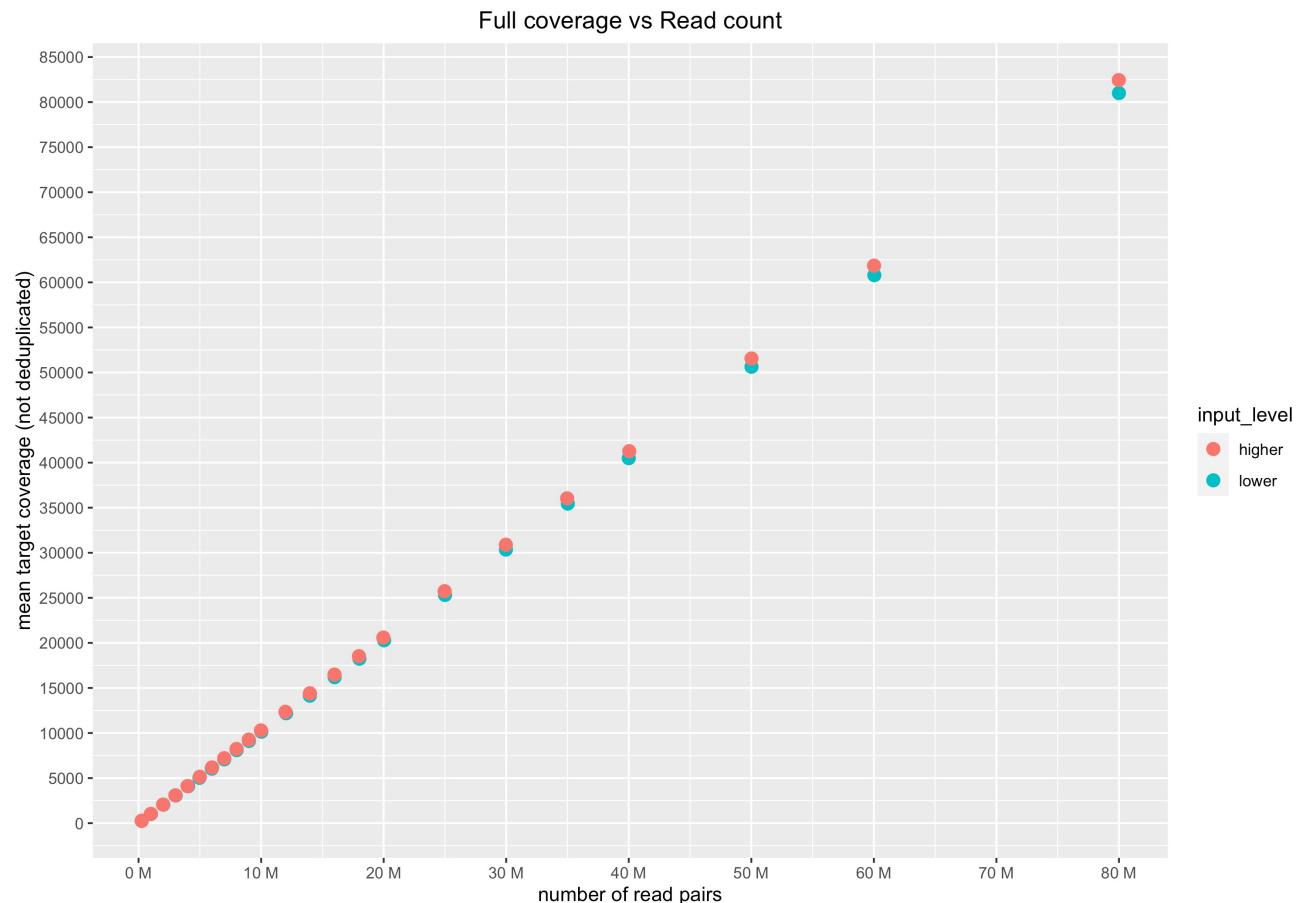
Duplication rate vs read count

- Increasing number of reads give increasing duplication rate
- Small difference in DNA input amount (few ng)
- Higher input amount gives lower duplication rate



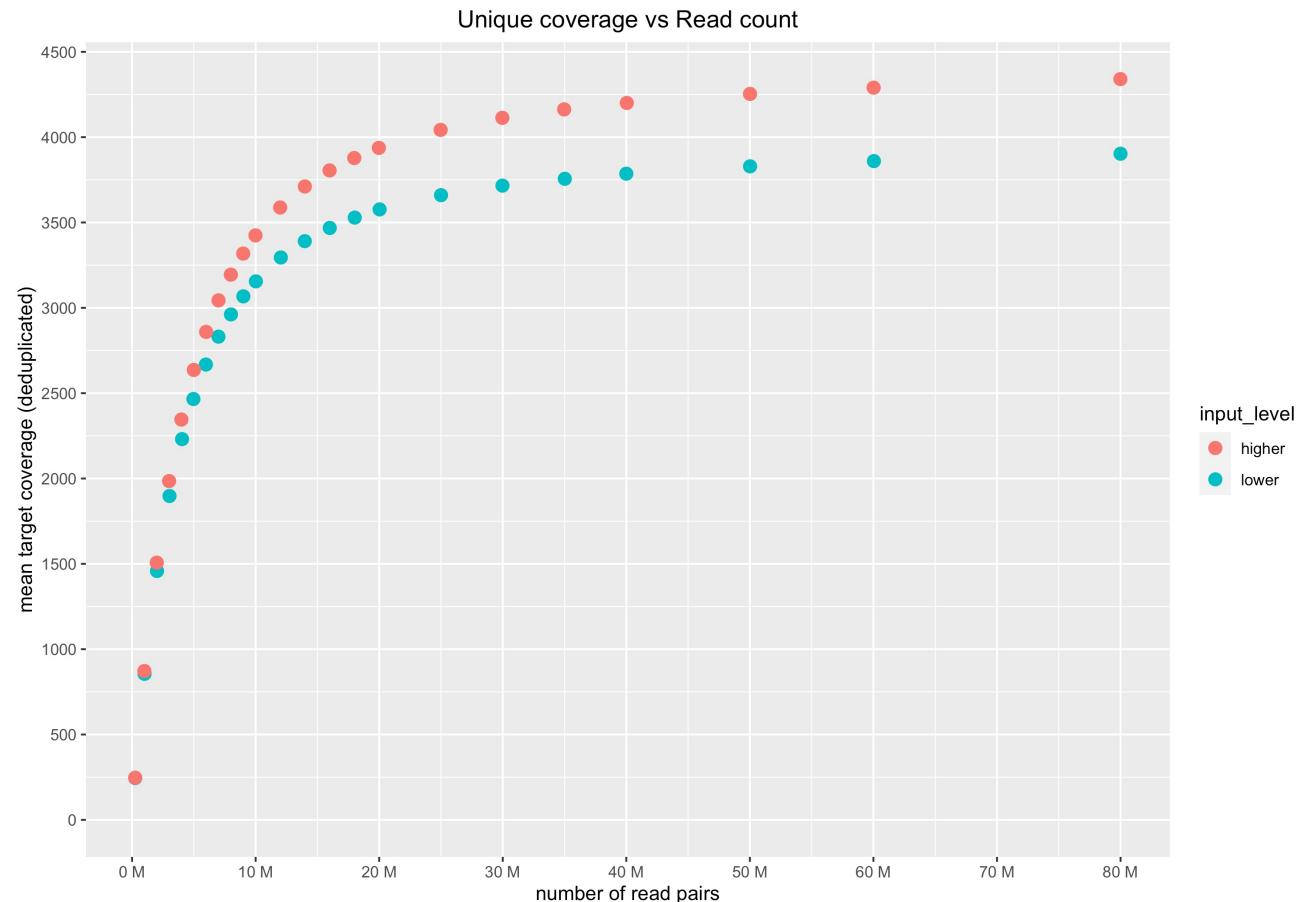
Coverage vs read count

- Including duplicates
 - Not de-duplicated
- Increasing number of reads give coverage
- No difference by input amount



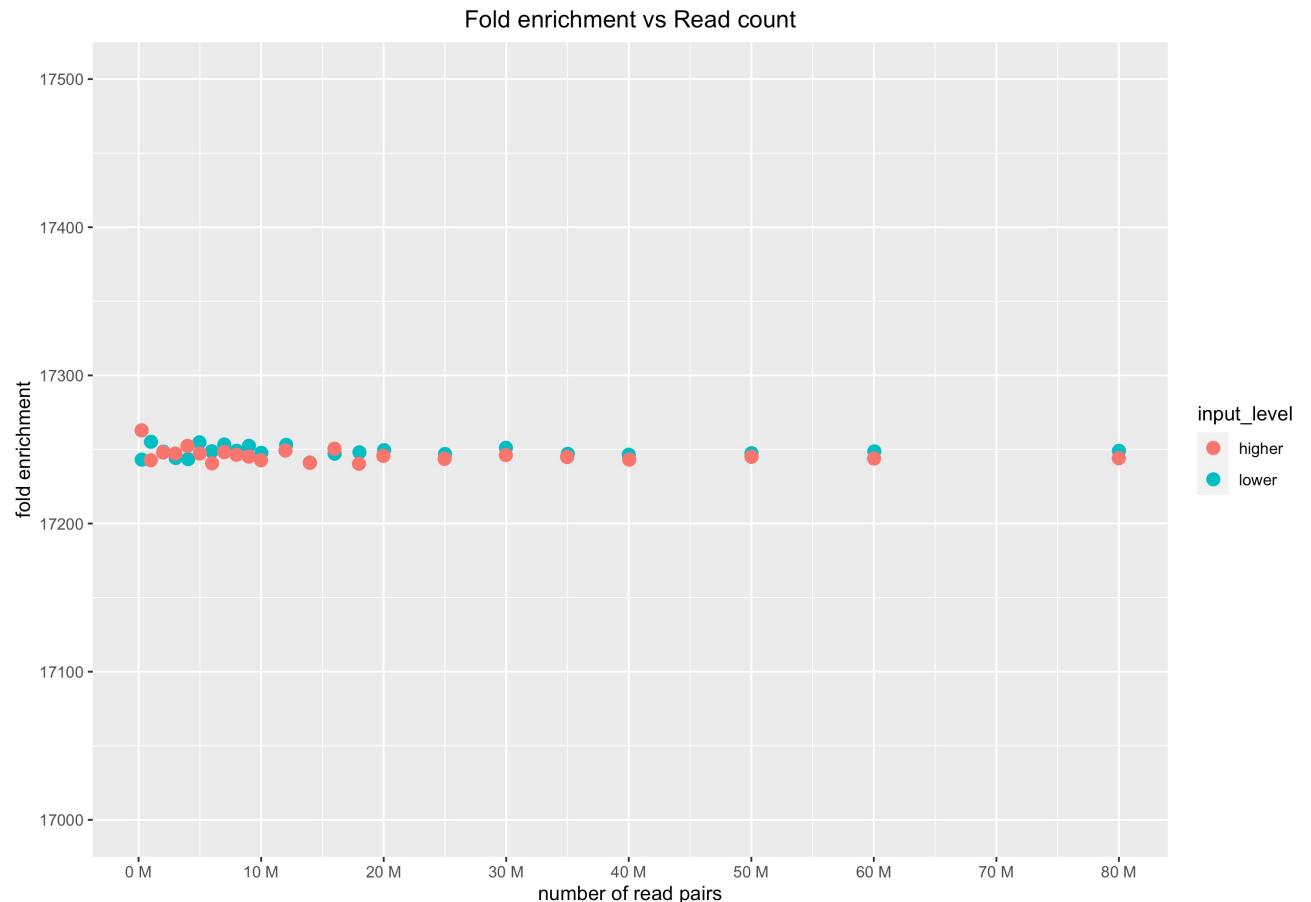
Coverage vs read count

- Unique
 - De-duplicated
- Increasing number of reads give coverage
- Higher input amount gives higher unique coverage
- How does the input amount relate to the max coverage possible?
- How does this relate to duplication rate?



Fold enrichment vs read count

- Before de-duplication
- Independent of read count and input amount
- Related to the size of targeted regions
- What could be the cause of a deviating number?



Quality control metrics – quiz

1. How does the DNA input amount relate to the max coverage possible?
 - a) Increased max coverage with decreased input amount
 - b) Decreased max coverage with decreased input amount
 - c) No correlation between them
2. How does this relate to duplication rate?
 - a) Increased duplication rate with decreased input amount
 - b) Decreased duplication rate with decreased input amount
 - c) No correlation between them
3. What could be the cause of a deviating fold enrichment?
 - a) Too low DNA input amount
 - b) Failed capture
 - c) Failure on the sequencer
4. How can these metrics help in answering “is it good enough”?
 - a) Was the min required coverage reached?
 - b) Did capture work, so we got data for the targeted regions?
 - c) Was the DNA input amount high enough?

Quality control metrics – quiz answers

1. How does the input amount relate to the max coverage possible?
 - b) Lower input → fewer unique molecules → lower max unique coverage
2. How does this relate to duplication rate?
 - a) Lower input → fewer unique molecules → higher duplication rate for same number of reads
3. What could be the cause of a deviating fold enrichment?
 - b) Failure in capture step
 - Lower fold enrichment than expected → more reads in non-targeted regions
4. How can these metrics help in answering “is it good enough”?
 - a) Certain min coverage required, if less → not good
 - b) Much lower fold enrichment → capture failure → not good
 - c) Low coverage and low duplication rate (~ <60%) → not all unique molecules sequenced → sequencing more will give higher coverage

QC metrics – Questions?

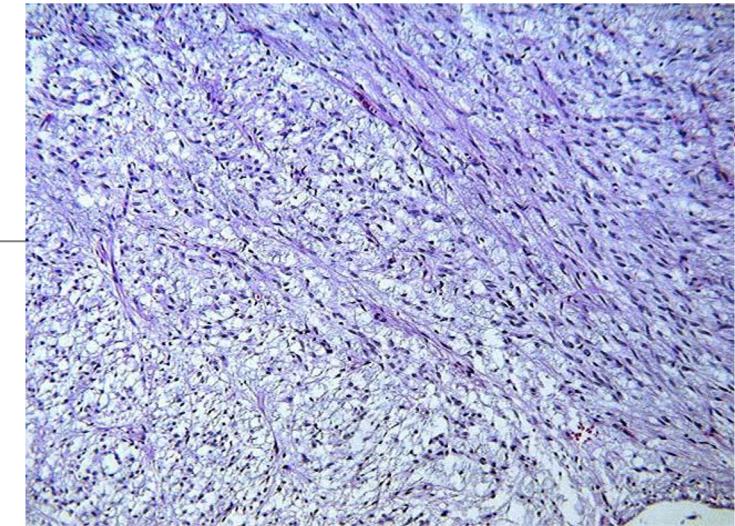


Karolinska
Institutet

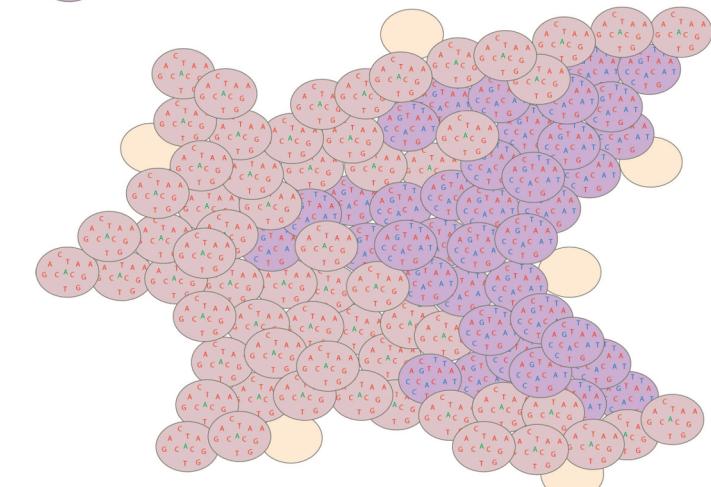
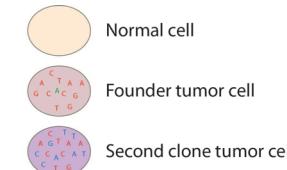
Somatic and germline variant calling

Tumor purity and clonality

- Tumors are often impure
 - Mix of tumor cells and normal (germline) cells
 - Tumor purity expressed as fraction or percentage
 - Cancer DNA fraction is a very closely related term
- Tumors often contain multiple clones
 - Diverse collections of cells harboring different mutations
 - Often one original clone with initial mutations
 - Subclones containing additional mutations may form
 - Treatments may favor one subclone which has resistance mutation, causing it to take over
 - → clonal evolution

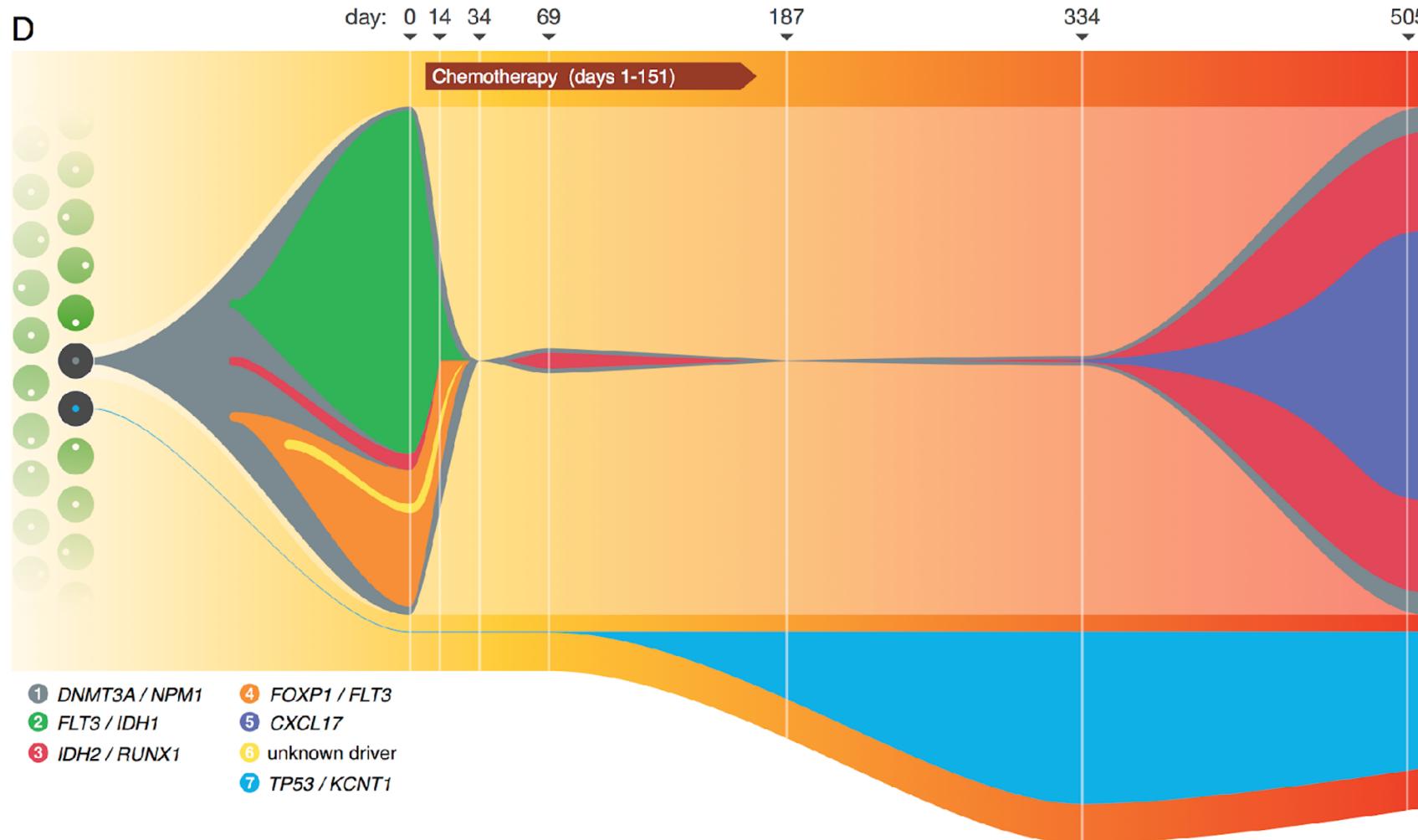


Credit Griffith brothers' lecture #5



Credit Griffith brothers' lecture #5

Clonal evolution



Somatic mutation vs. germline mutation vs. germline polymorphism

- Germline mutations
 - Present in egg or sperm
 - All cells of affected offspring
 - Heritable
 - Cause of familial cancers
- Germline polymorphisms
 - Present in egg or sperm
 - All cells of affected offspring
 - Heritable
 - Common in the population, > 1%
 - Generally not impacting disease
 - SNP – Single Nucleotide Polymorphism
- Somatic mutations
 - Occur in non-germline tissues
 - Only tumor cells (breast, lung, blood, etc.)
 - Non-heritable
 - Various reasons
 - Smoking, UV light, oxidation in cells, etc.
 - Cause of sporadic cancers
 - And of familial cancers – in combination with germline mutations

Small variants

- SNVs – single nucleotide variants
 - Change from one base to another at one single position of the genome
- Indels – Small insertions and deletions
 - 1 – ~30 bases added or removed

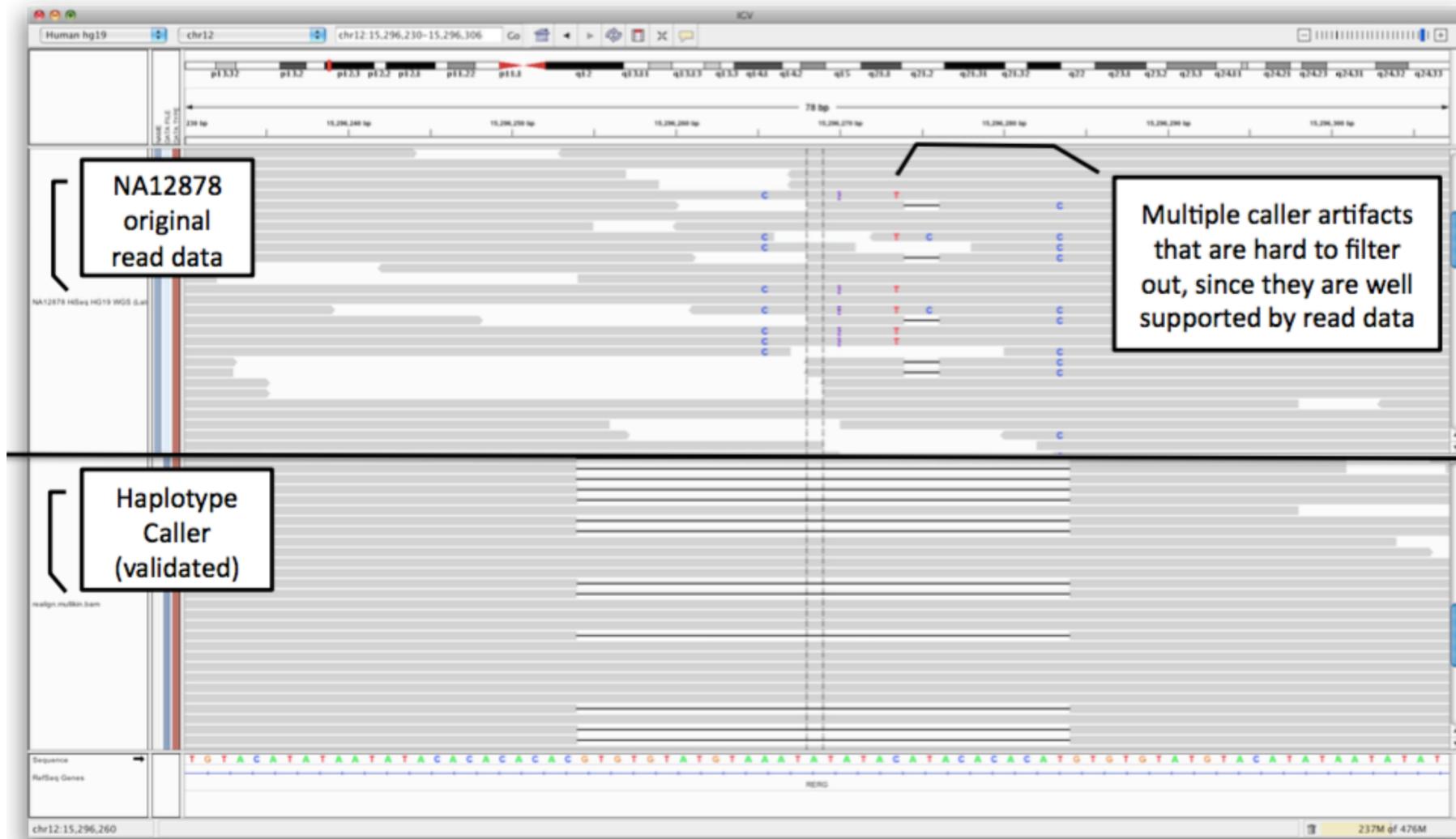
Germline mutation calling – methodology

- Mutation/variant calling – to identify mutations from e.g. sequencing data
 - Focus: DNA data
- Germline sample: often white blood cells, WBCs
- SNVs and small indels
- VAFs (variant allele frequencies): ~50% or 100%
 - One or two mutated alleles out of totally two alleles
- General method:
 - Find positions where a fraction of mapped reads have base deviating from the reference genome – alternate allele
 - If significant difference from the reference – call a germline variant

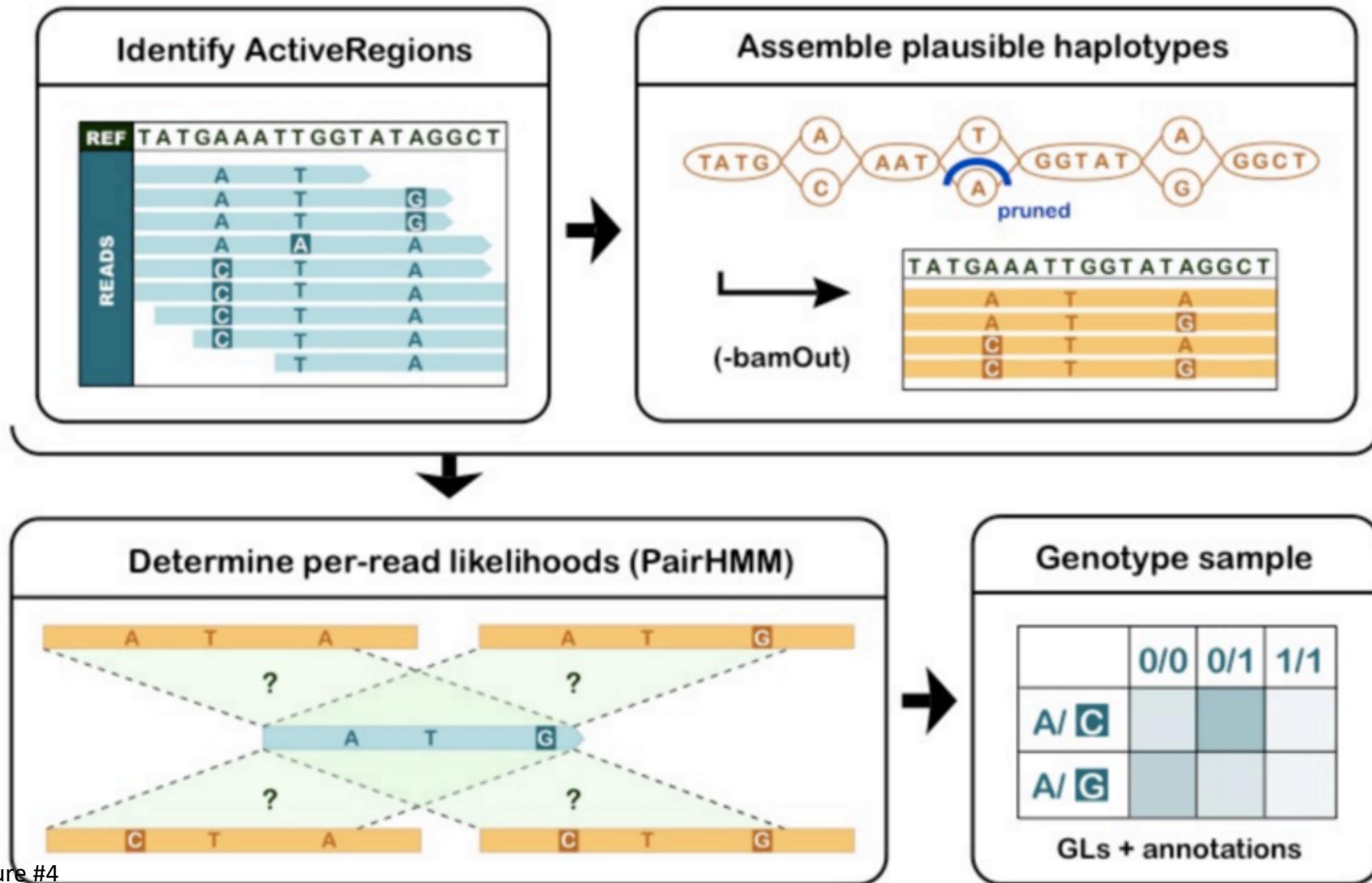
Germline mutation calling – methodology and tools

- Call SNVs and Indels separately by considering each variant locus
 - Very fast
 - Assumes bases are independent
- Call SNVs and indels simultaneously via Bayesian genotype likelihood model
 - More computationally intensive
- Call SNVs, indels and SVs simultaneously by performing a local de novo assembly
 - More computationally intensive
 - More accurate – gets rid of many false positives especially indels
 - GATK HaplotypeCaller
 - Strelka

Germline mutation calling – methodology



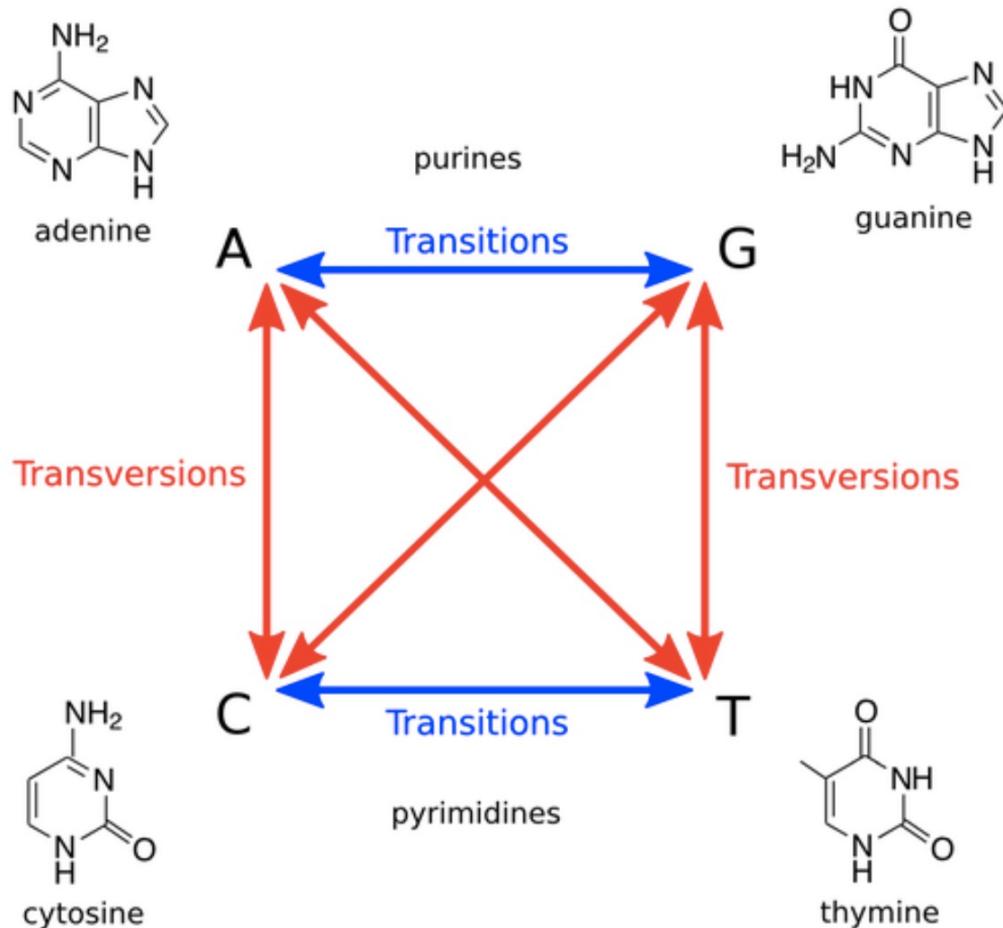
Germline mutation calling – HaplotypeCaller



GATK recommended filters

- SNPs
 - QD < 2.0 (variant quality/depth of non-ref samples)
 - MQ < 40.0 (Mapping quality)
 - FS > 60.0 (Phred score Fisher's test pvalue for strand bias)
 - SOR > 3.0 (Strand odds ratio, aims to evaluate whether there is strand bias in the data)
 - MQRankSum < -12.5 (mapping quality of reference reads vs alt reads)
 - ReadPosRankSum < -8.0 (distance of alt reads from end of the read)
- Indels
 - QD < 2.0
 - FS > 200.0
 - SOR > 10.0
 - ReadPosRankSum < -20.0
 - InbreedingCoeff < -0.8

Transition/Transversion ratio (Ti/Tv)



Ratios:

Random = 0.5

WGS = 2.0-2.1

Exome = 3-3.5

Watch for major deviation from typical ratio

A good paper

- The “gnomAD” paper
- Published in Nature
- Supplementary Information:
Many details on how to properly call and filter germline variants

Article

The mutational constraint spectrum quantified from variation in 141,456 humans

<https://doi.org/10.1038/s41586-020-2308-7>

Received: 27 January 2019

Accepted: 26 March 2020

Published online: 27 May 2020

Open access

 Check for updates

Konrad J. Karczewski^{1,2}✉, Laurent C. Francioli^{1,2}, Grace Tiao^{1,2}, Beryl B. Cummings^{1,2,3}, Jessica Alföldi^{1,2}, Qingbo Wang^{1,2,4}, Ryan L. Collins^{1,4,5}, Kristen M. Laricchia^{1,2}, Andrea Ganna^{1,2,6}, Daniel P. Birnbaum^{1,2}, Laura D. Gauthier⁷, Harrison Brand^{1,5}, Matthew Solomonson^{1,2}, Nicholas A. Watts^{1,2}, Daniel Rhodes⁸, Moriel Singer-Berk^{1,2}, Eleina M. England^{1,2}, Eleanor G. Seaby^{1,2}, Jack A. Kosmicki^{1,2,4}, Raymond K. Walters^{1,2,9}, Katherine Tashman^{1,2,9}, Yossi Farjoun⁷, Eric Banks⁷, Timothy Poterba^{1,2,9}, Arcturus Wang^{1,2,9}, Cotton Seed^{1,2,9}, Nicola Whiffin^{1,2,10,11}, Jessica X. Chong¹², Kaitlin E. Samocha¹³, Emma Pierce-Hoffman^{1,2}, Zachary Zappala^{1,2,14}, Anne H. O'Donnell-Luria^{1,2,15,16}, Eric Vallabh Minikel¹, Ben Weisburd⁷, Monkol Lek¹⁷, James S. Ware^{1,10,11}, Christopher Vittal^{2,9}, Irina M. Armean^{1,2}, Louis Bergelson⁷, Kristian Cibulskis⁷, Kristen M. Connolly¹⁸, Miguel Covarrubias⁷, Stacey Donnelly¹, Steven Ferriera¹⁸, Stacey Gabriel¹⁸, Jeff Gentry⁷, Namrata Gupta¹⁸, Thibault Jeandet⁷, Diane Kaplan⁷, Christopher Llanwarne⁷, Ruchi Munshi⁷, Sam Novod⁷, Nikelle Petrillo⁷, David Roazen⁷, Valentin Ruano-Rubio⁷, Andrea Saltzman¹, Molly Schleicher¹, Jose Soto⁷, Kathleen Tibbetts⁷, Charlotte Tolonen⁷, Gordon Wade⁷, Michael E. Talkowski^{1,5,19}, Genome Aggregation Database Consortium^{*}, Benjamin M. Neale^{1,2,9}, Mark J. Daly^{1,2,6,9} & Daniel G. MacArthur^{1,2,150,151}✉

Germline mutation calling – quiz

1. Are germline mutations heritable?
 - a) Yes
 - b) No
 - c) Only sometimes
2. Which of these variant allele frequencies is most likely for a germline variant in a germline sample?
 - a) 25%
 - b) 50%
 - c) 75%
3. What is a haplotype?
 - a) A germline variant caller
 - b) A very common germline variant
 - c) A possible combination of nearby germline variants

Germline mutation calling – quiz answers

1. Are germline mutations heritable?
 - a) Yes – these are the mutations we are born with, and they can be inherited by our children – they are in our germ cells
2. Which of these variant allele frequencies is most likely for a germline variant in a germline sample?
 - b) 50% - 1 mutated out of 2 alleles
3. What is a haplotype?
 - c) A possible combination of nearby germline variants
 - HaplotypeCaller is a germline variant caller that calls variants by testing the likelihood of possible haplotypes based on the DNA read support

Germline variant calling – Questions?

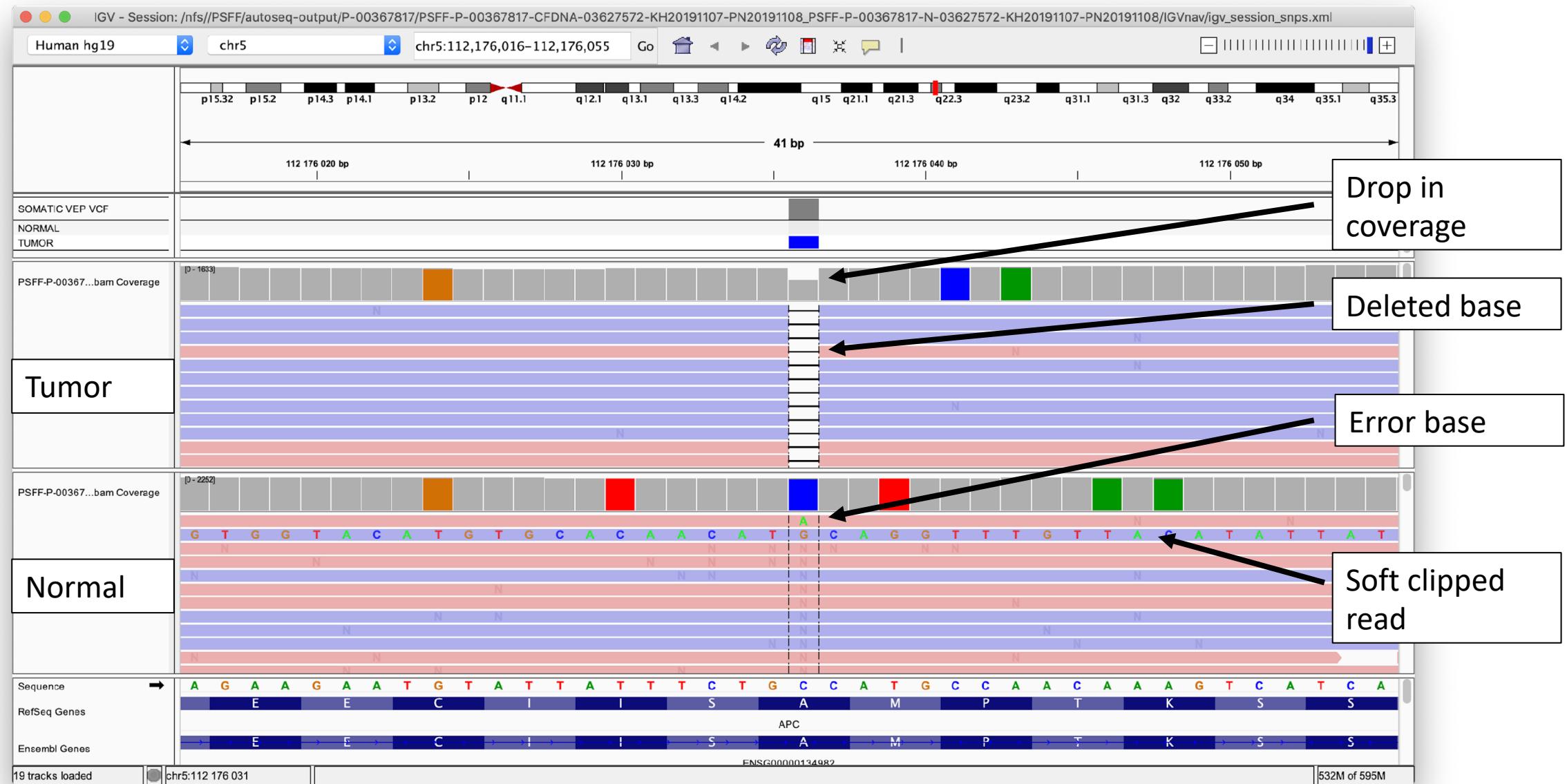
Somatic mutation calling – methodology

- Mutation/variant calling – to identify mutations from e.g. sequencing data
 - Focus: DNA data
- Somatic mutations are best distinguished by comparison of tumor to a matched normal
- Normal sample: Germline sample from the same individual
 - If not available: use healthy donor sample → requires additional filtering
- SNVs and small indels
- VAFs: ~1% - 100%, depending on purity, ploidy, clonality
 - Limited detection capacity for VAF < 1%
- General method:
 - Find positions where a fraction of mapped reads have base deviating from the reference genome – alternate allele
 - Compare with germline sample
 - If significant difference to germline – call a somatic variant

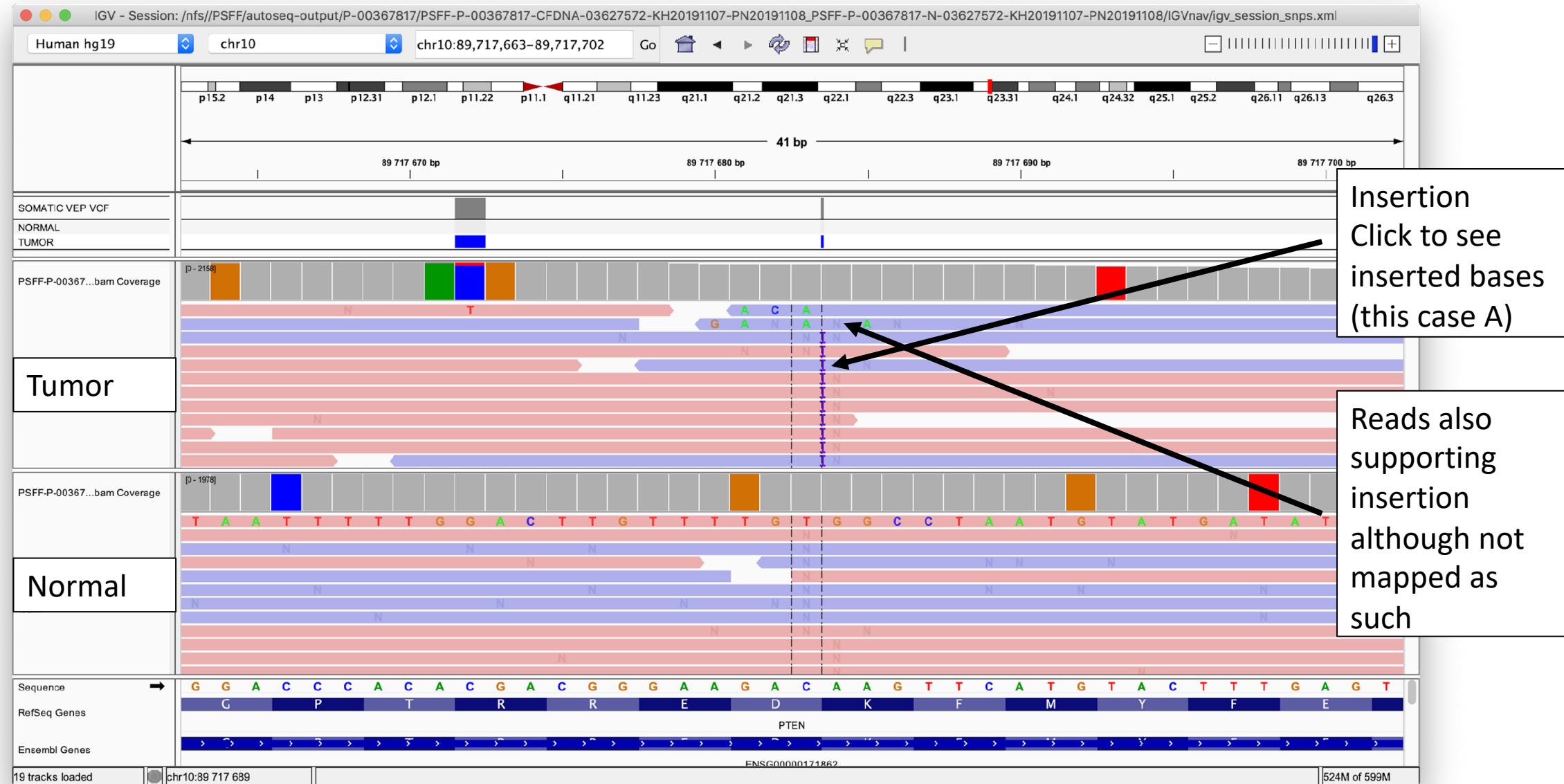
Somatic mutation calling – SNV



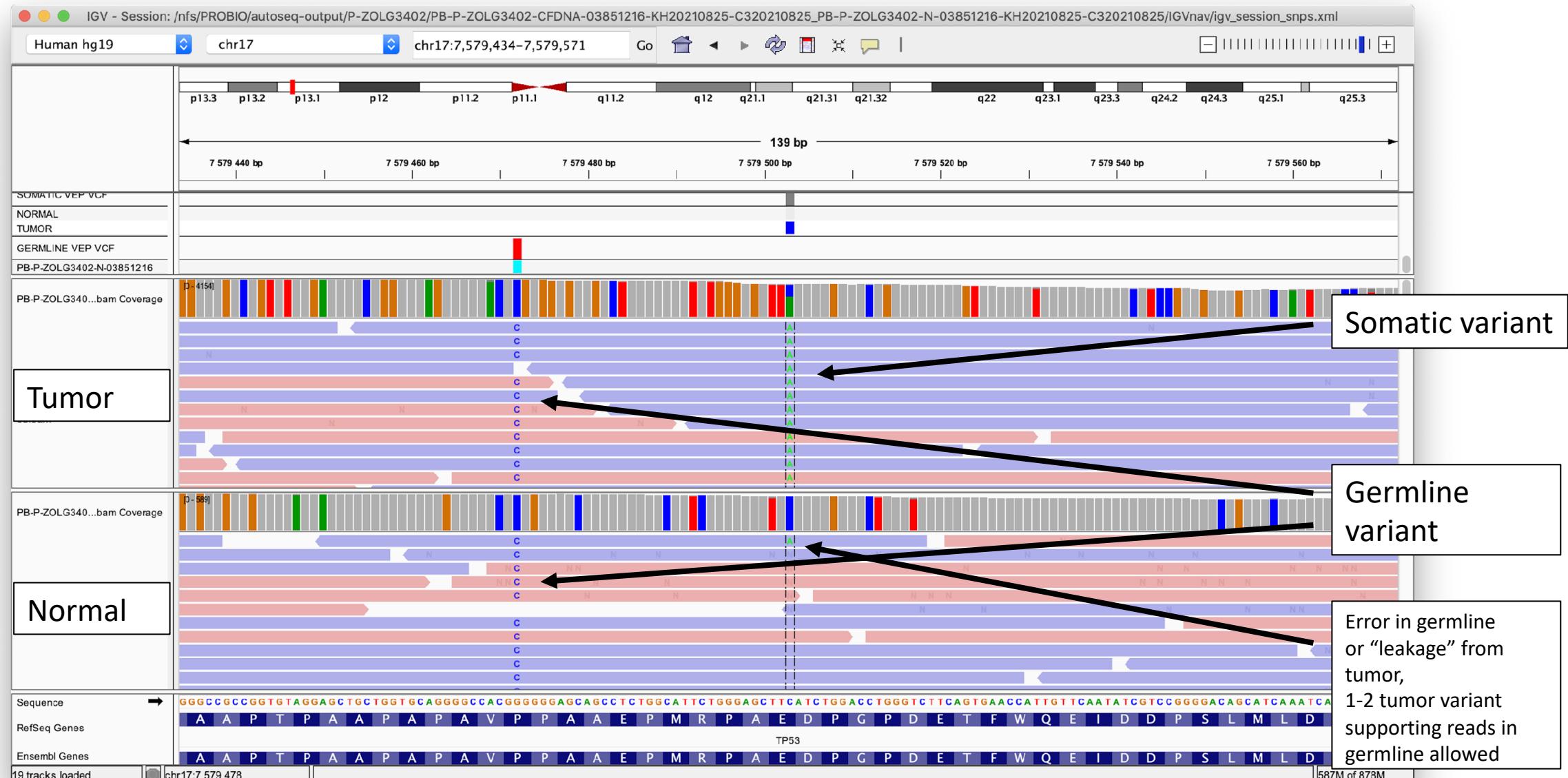
Somatic mutation calling – deletion



Somatic mutation calling – insertion



Germline vs somatic variant



Somatic mutation calling – tools

- Tools (variant callers):
 - Vardict, Varscan – counting the reads with ref bases and non-ref bases (alt allele) in tumor and normal, calculating significance test, if significant call the variant
 - Strelka, Mutect2 – Local de novo assembly
 - Different properties – better at different types of variants

VarScan 2: Somatic mutation and copy number alteration discovery in cancer by exome sequencing

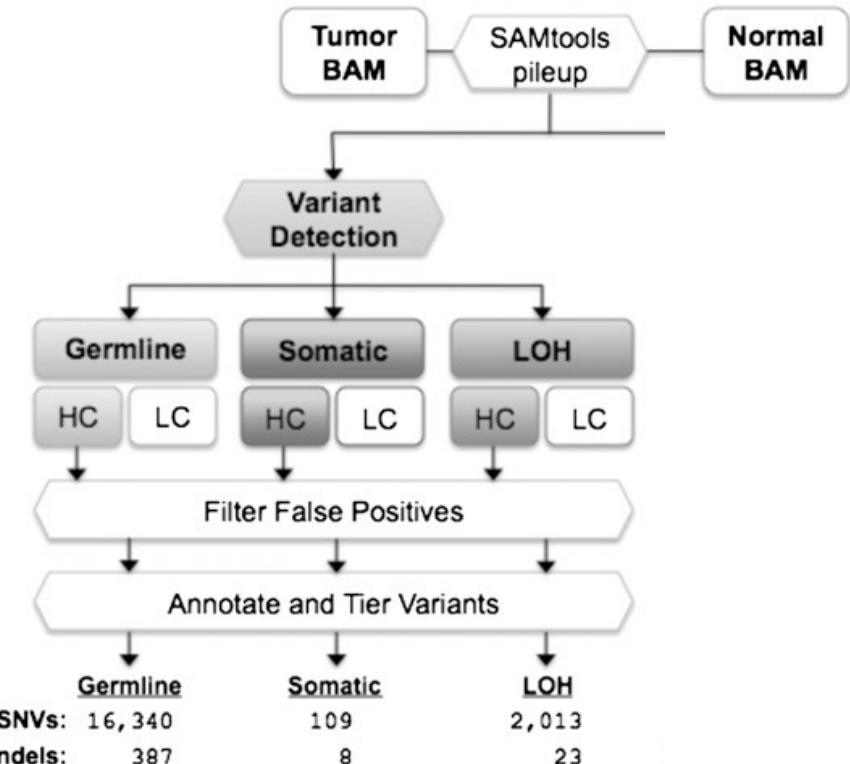
Daniel C. Koboldt,¹ Qunyuan Zhang,¹ David E. Larson,¹ Dong Shen,¹
Michael D. McLellan,¹ Ling Lin,¹ Christopher A. Miller,¹ Elaine R. Mardis,^{1,2,3} Li Ding,^{1,2,4}
and Richard K. Wilson^{1,2,3,4}

¹*The Genome Institute, Washington University, St. Louis, Missouri 63108, USA;* ²*Department of Genetics, Washington University, St. Louis, Missouri 63110, USA;* ³*Siteman Cancer Center, Washington University, St. Louis, Missouri 63110, USA*

Published in Genome Research

VarScan 2

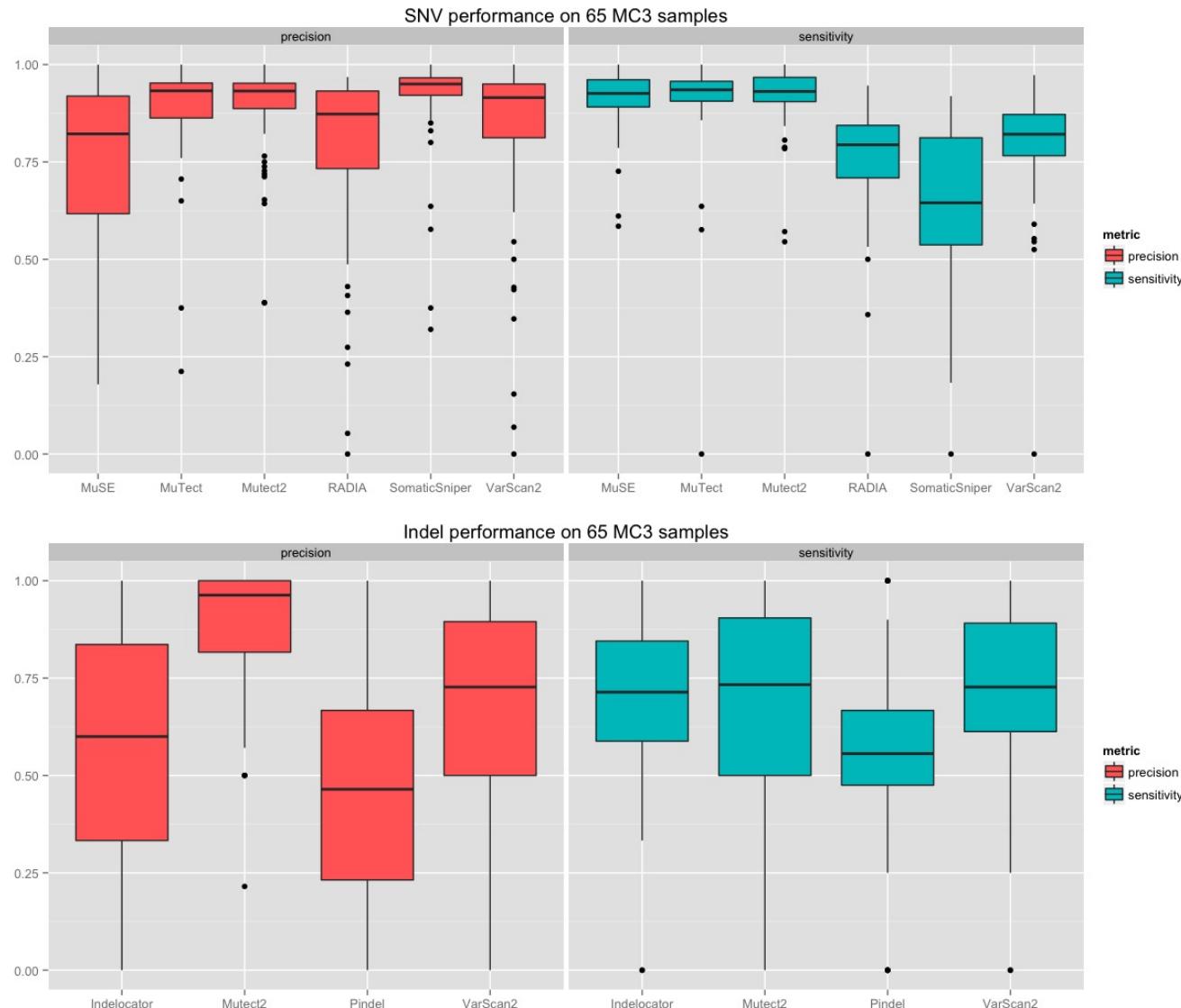
- Samtools pileup
 - Checks every position in given regions, which base every read covering it has (A, C, G, T, ins, del) and the total depth
- Variant detection:
 - Thresholds for coverage, base quality, variant allele frequency
 - Statistically significant variant: Fisher's exact test of ref and alt read counts, compared to expected sequencing error distribution
- Variant classification:
 - Comparison between variants in germline and in tumor
 - Fisher's exact test of ref and alt read counts in tumor and in normal
 - If variant only in tumor → somatic variant
 - If heterozygous ($VAF < 75\%$) in germline but homozygous in tumor ($VAF > 75\%$) → LOH variant (loss of heterozygosity, germline variant where the ref allele is somehow lost in the tumor)
 - If variant has same genotype (heterozygous/homozygous) in both tumor and germline → germline variant
 - Filters on read position, strandedness, variant reads, variant frequency, distance to 3', homopolymer, mapping quality difference, read length difference, MMQS (mismatch quality sum) difference



Fisher's test 2x2	Germline sample	Tumor sample
Reference allele	572	585
Alternate allele	1	213

Comparison of callers

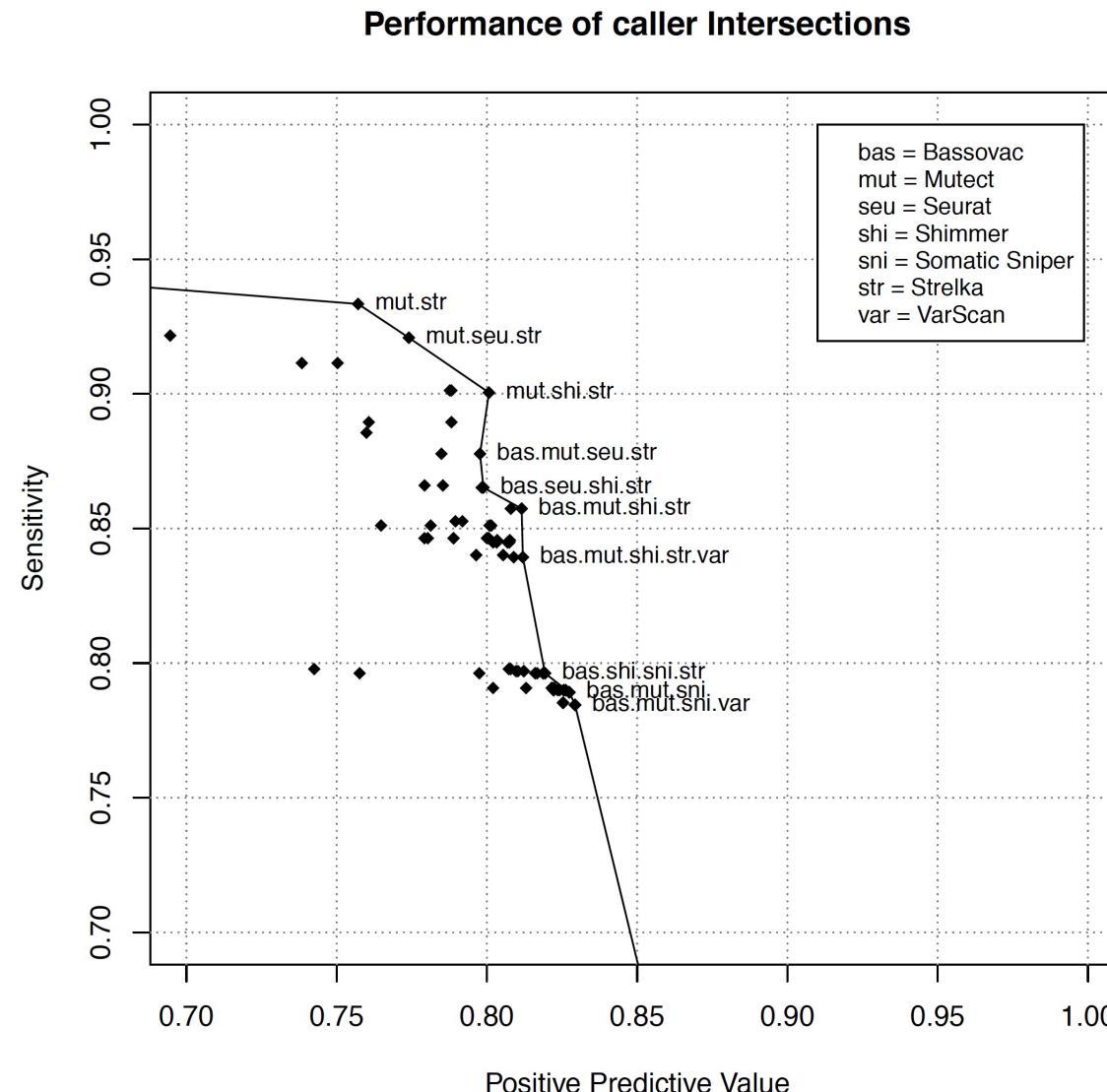
- Callers have different performance on different variant types
- Comparison from Benjamin, D., Sato, T., Cibulskis, K., Getz, G., Stewart, C., & Lichtenstein, L. (2019). Calling Somatic SNVs and Indels with Mutect2. *BioRxiv*. <https://doi.org/10.1101/861054>



Somatic mutation calling – tools

- Consensus:
 - Not all called variants are true
 - Take consensus from the four callers to improve specificity
 - Literature review in 2018 showed that these (Mutect, Strelka, Varscan, Vardict) were the most commonly used and give a good combination of properties
 - Tool somaticseq
 - If ≥ 2 callers call the same mutation, include it in table for manual curation

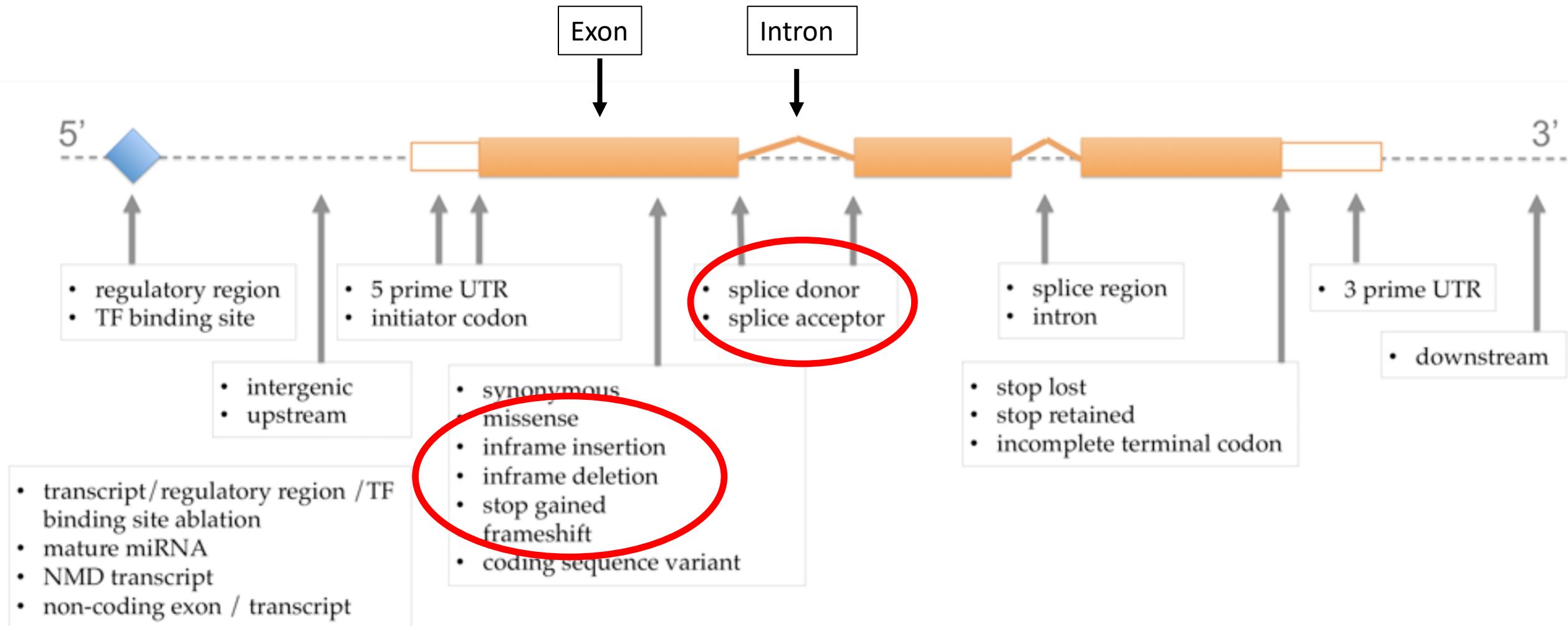
Somatic mutation calling – tools



Somatic mutation calling – annotation

- What effects do the mutations have?
 - On proteins
 - On disease
 - On treatment alternatives
- VEP – Variant Effect Predictor
 - From Ensembl
 - “VEP determines the effect of your variants ... on genes, transcripts, and **protein sequence**, as well as regulatory regions.”
(<https://www.ensembl.org/info/docs/tools/vep/index.html>)

Somatic mutation calling – VEP annotation categories



Mutation calling – file format

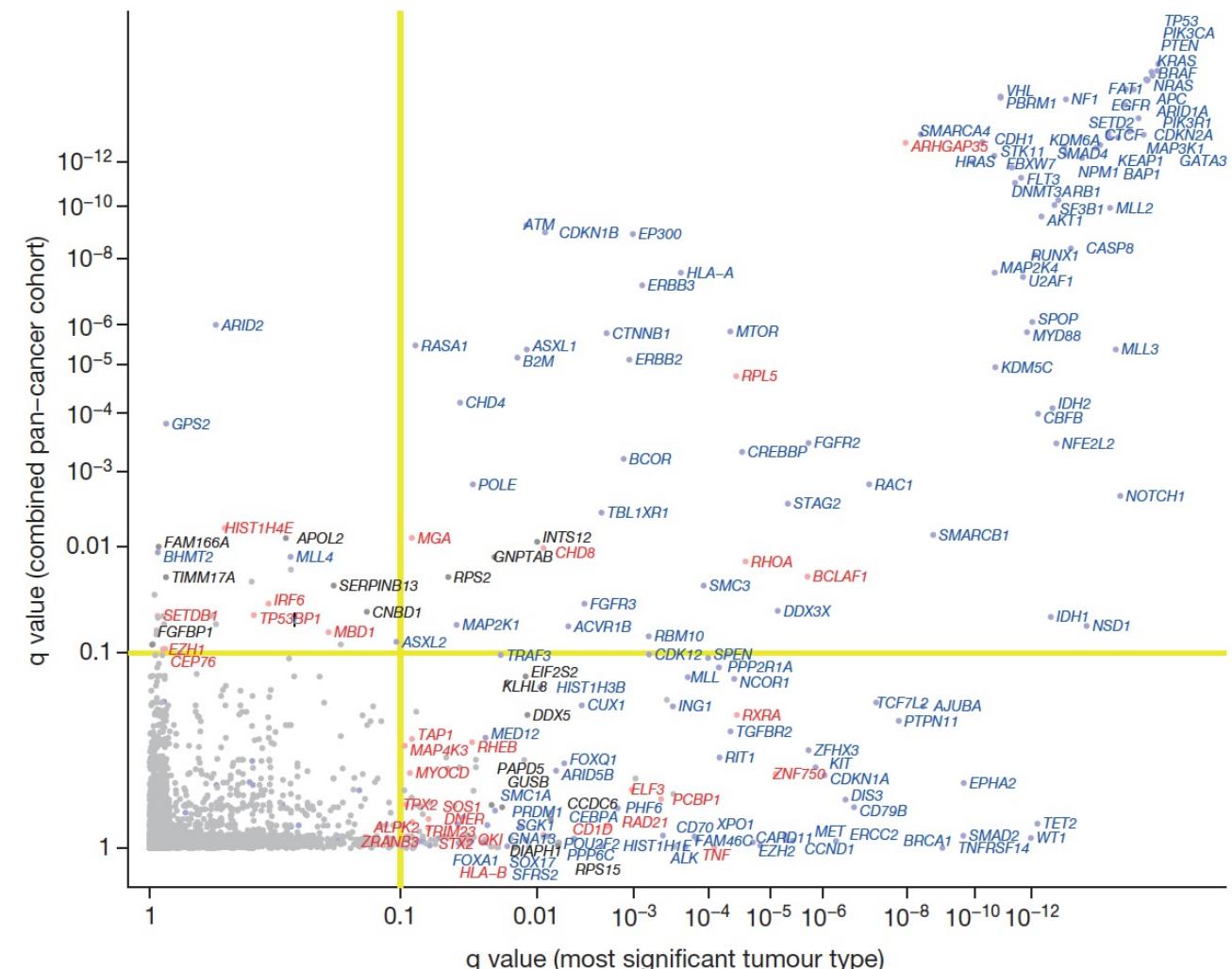
- VCF – variant call format
 - Header – metadata describing different fields
 - Main

Mutation calling – manual curation

- Not all called variants are true – even after consensus
- If filters were sharpened to decrease number of false positives, we would instead miss some variants
- Some properties are hard to account for/discover programmatically
- Manual curation necessary
- Purpose: to identify true variants with impact on the protein structure

Different genes mutated in different cancers

- x-axis: how common in the tumor type where it's most common?
 - y-axis: how common in all tumor types together?
 - Above yellow line → candidate cancer gene
 - Many genes that are candidate cancer genes when looking at a specific tumor type are not significantly mutated when looking at all cancer types together
 - Lawrence, M. S., Stojanov, P., Mermel, C. H., Robinson, J. T., Garraway, L. A., Golub, T. R., Meyerson, M., Gabriel, S. B., Lander, E. S., & Getz, G. (2014). Discovery and saturation analysis of cancer genes across 21 tumour types. *Nature*, 505(7484), 495–501. <https://doi.org/10.1038/nature12912>



Somatic mutation calling – quiz

1. Are somatic mutations heritable?
 - a) Yes
 - b) No
 - c) Only sometimes
2. In which sample can you find the somatic mutations?
 - a) Tumor sample
 - b) Germline sample
 - c) Both
3. What is a frameshift variant?
 - a) A single position base change, that gives a different codon and amino acid
 - b) Insertion of a new stop codon before the actual stop of the gene, causing premature stop of the protein
 - c) Insertion or deletion that changes the open reading frame, and thus all downstream amino acids
4. What is a VCF?
 - a) A file listing the regions of interest, used as input to variant caller
 - b) A file containing variants, e.g. as output of variant caller
 - c) A file with positions that don't have any variant, so called wild type positions

Somatic mutation calling – quiz answers

1. Are somatic mutations heritable?
 - b) No
2. In which sample can you find the somatic mutations?
 - a) Tumor sample
3. What is a frameshift variant?
 - c) Insertion or deletion that changes the open reading frame, and thus all downstream amino acids
 - One codon is 3 bases, if indel size is not multiple of 3 the downstream codons won't be read correctly, but under a different reading frame
4. What is a VCF?
 - b) A file containing variants, e.g. as output of variant caller



Karolinska
Institutet

Somatic variant calling – Questions?

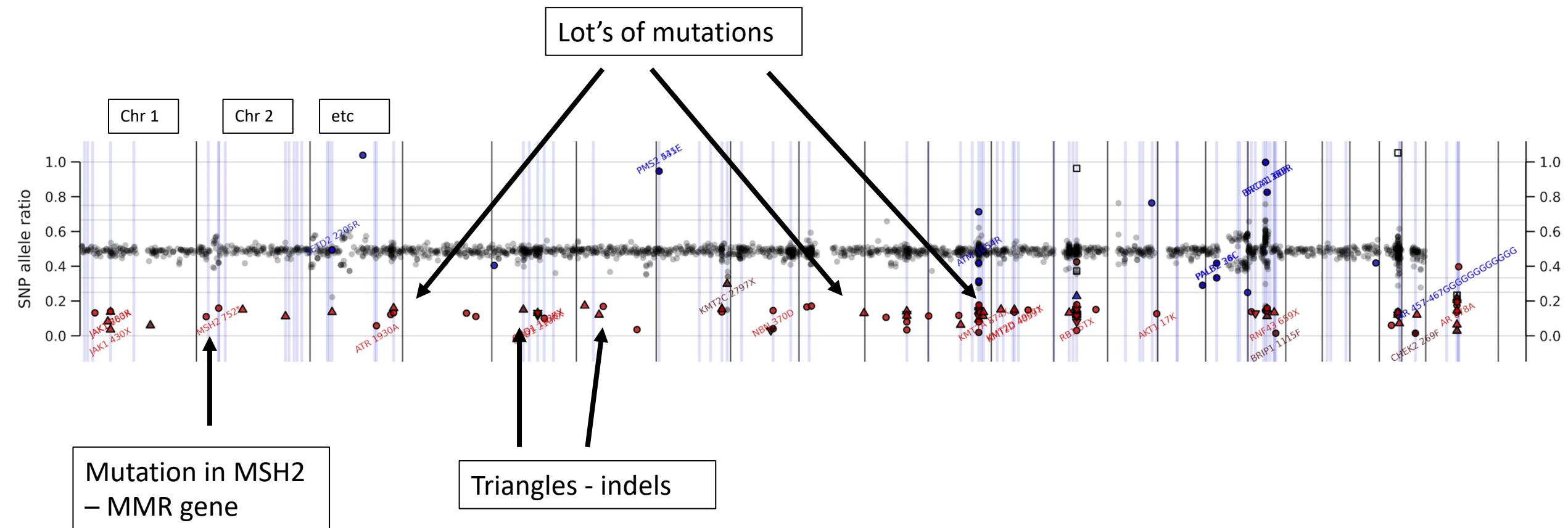
CHIP - Clonal Hematopoiesis of Indeterminate Potential

- Sub-population of blood cells carrying the same mutation(s)
- Age-related
- Increased risk of blood cancer and cardiovascular disease
- Shows in germline DNA, often at VAFs \neq 50% or 100%
- May show in cfDNA but not tissue tumor samples

MSI – microsatellite instability

- Microsatellite: small repetitive sequence of the genome, e.g. TTTTT or TATATATA
- MSI: increased rate of insertion or deletion of repeated segments in tumor
 - → plenty of insertions and deletions
- Caused by dysfunctional mismatch repair (MMR) mechanism
 - Indels of repeated segments not corrected
- Levels:
 - MSI-H: microsatellite instability high
 - (MSI-L: microsatellite instability low)
 - MSS: microsatellite stable
- Many mutations → tumor cells may express lots of weird proteins on their surfaces
 - Immunotherapy may be effective

MSI – an example



MSI - tool

- mSINGS - MicroSatellite Instability detection by Next Generation Sequencing
- Given list of microsatellites, 63 in our designs
- Background control created by ~20 healthy donor samples
- Comparing number of repeats in each microsatellite locus in tumor sample to the same sites in background control
- Locus called as unstable if significantly different from control samples
- Fraction of unstable loci gives mSINGS score
 - Threshold for MSI-H: > ~0.2
 - Confirm by visual inspection of mutation plot

Hypermutation

- MSI is a type of hypermutation – phenotypes with highly increased levels of mutation frequency
- Another type is caused by defective DNA replication repair (mismatch repair), due to mutations in DNA polymerases
- This gives a different mutational signature, with more SNVs than indels
- Other types of hypermutation can be caused by environmental factors (e.g. UV light, smoking, chemotherapy), and are associated with specific mutational signatures

Credits

- Malachi Griffith, Obi Griffith, Zachary Skidmore, Huiming Xia
 - Lecture notes from the course “Introduction to bioinformatics for DNA and RNA sequence analysis (IBDR01)”, 29 October – 2 November, 2018
 - McDonell Genome Institute, Washington University of St Louis School of Medicine

More questions?