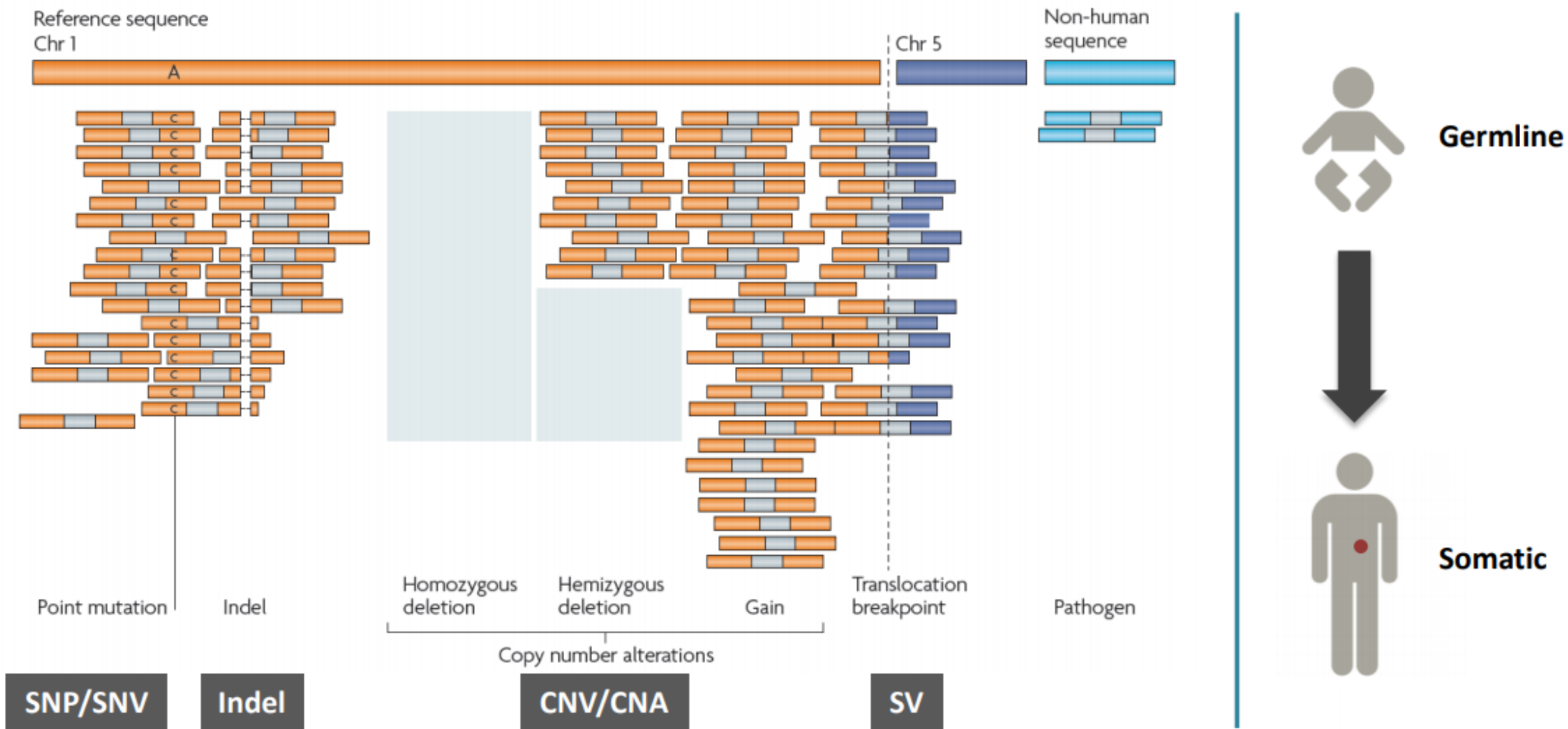


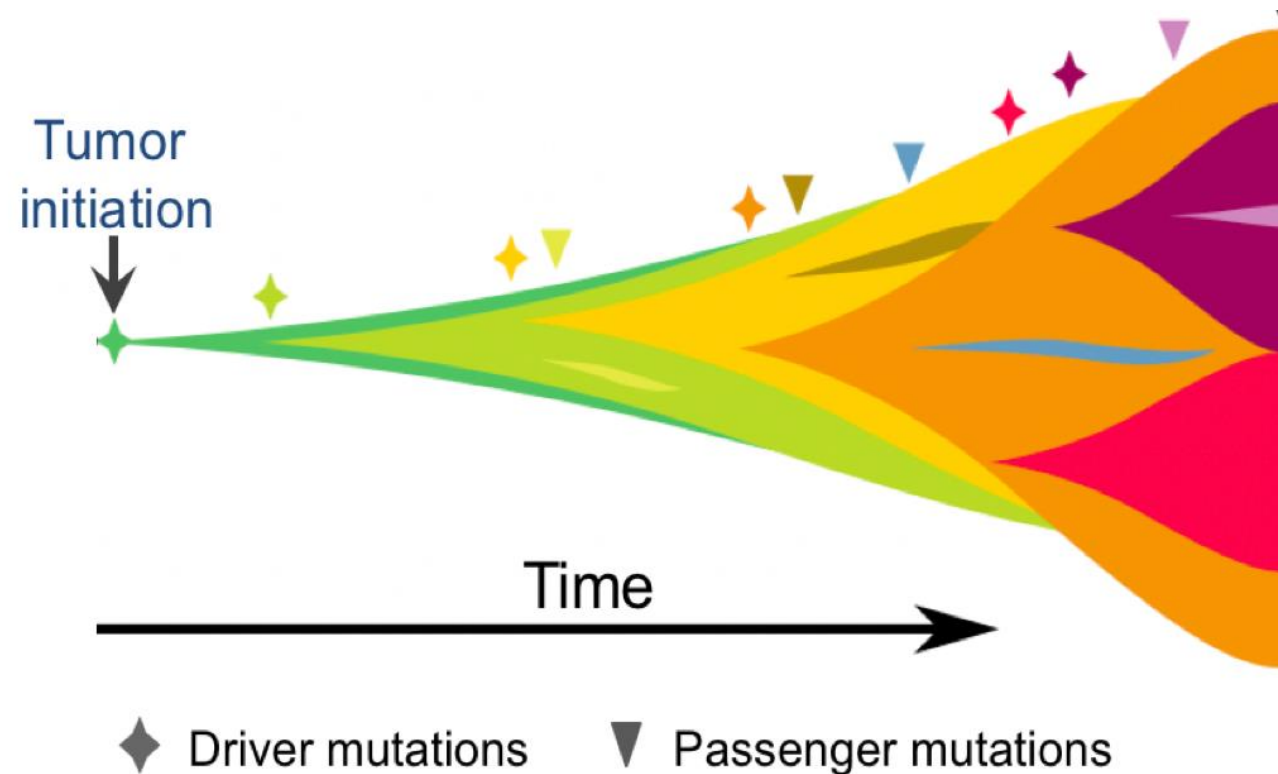
Somatic Variant Calling

Different types of variants

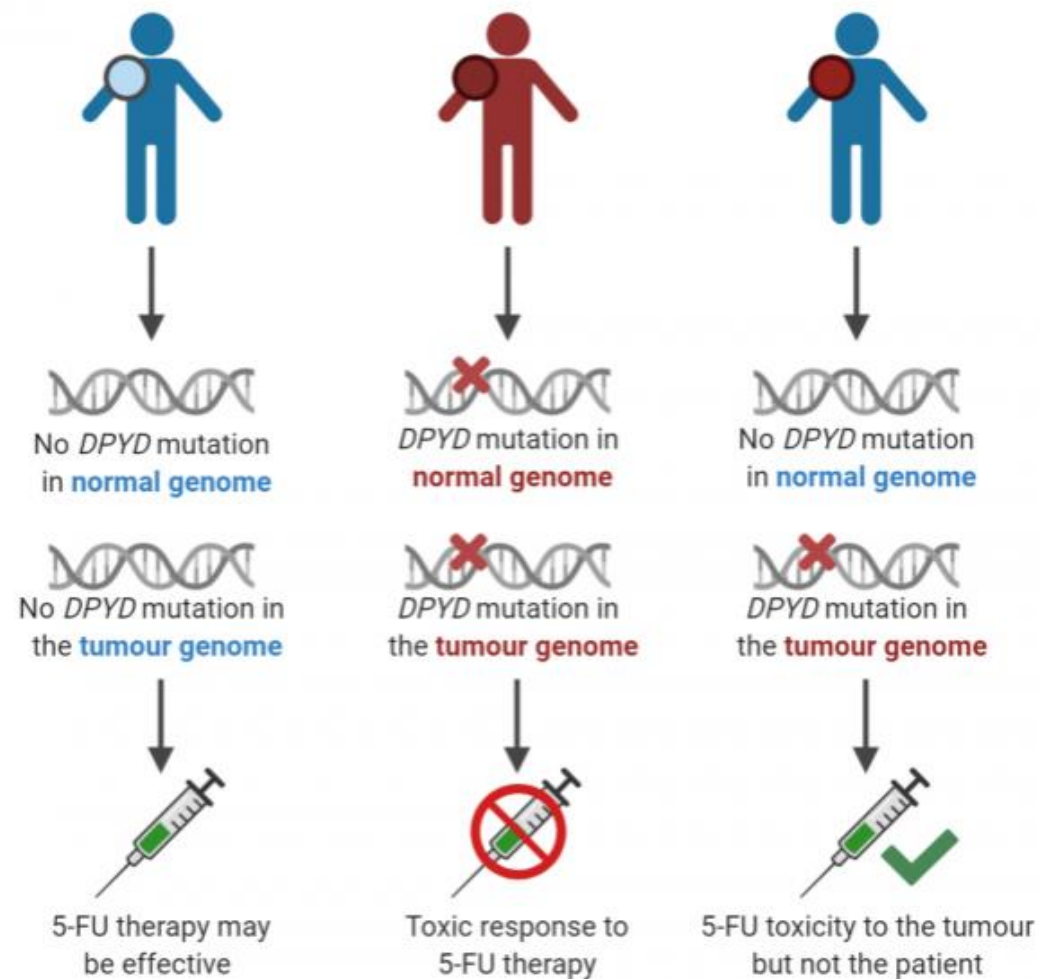
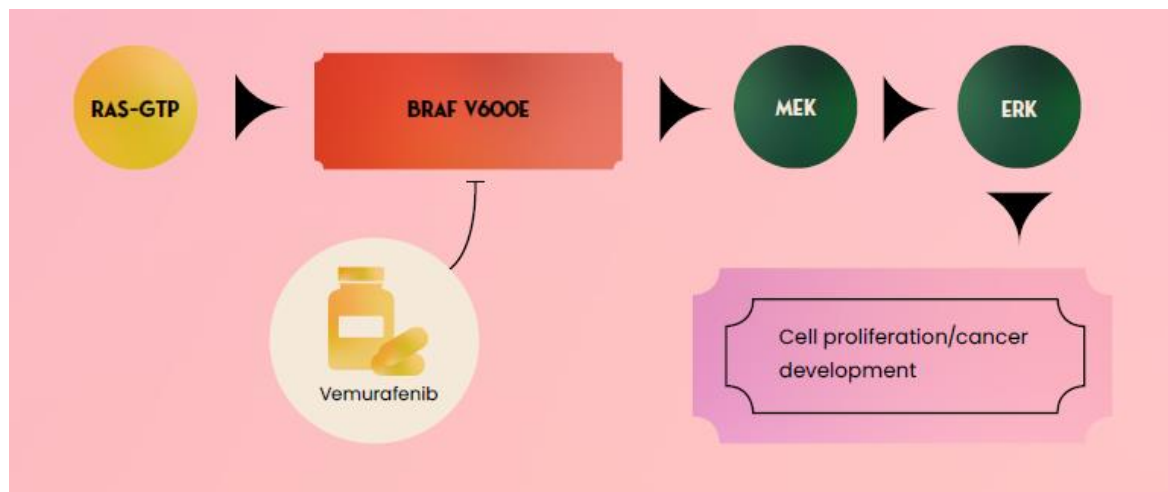


Role of somatic mutation events in tumor progression

Increasing tumor heterogeneity as genomic instability increases



Cancer sequencing helps guide and prioritize cancer treatment options



Why is SNV calling in cancer samples more complicated ?

FOR MANY REASONS

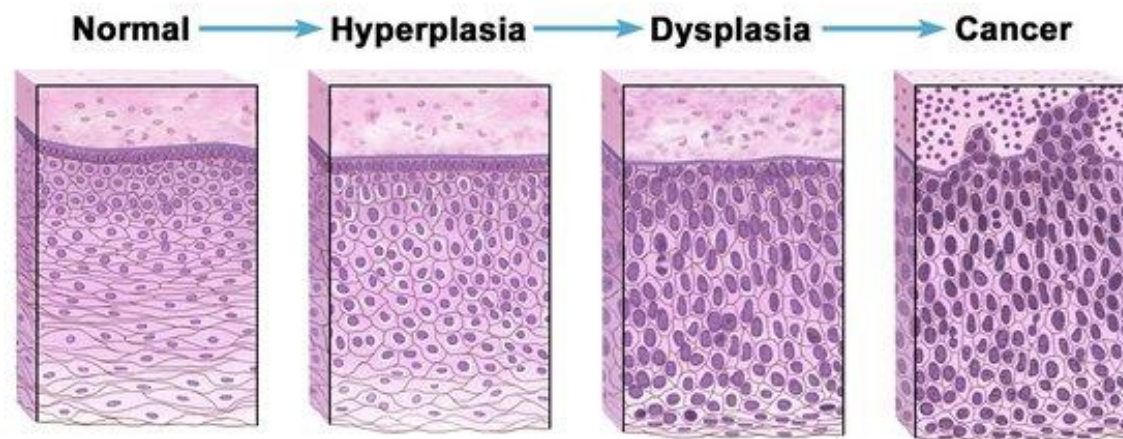
Why is SNV calling in cancer samples more complicated ?

Low tumor cellularity (tumour DNA content)

Tumor samples may have lower DNA content

→ need sensitivity in variant calling +++

Normal Cells May Become Cancer Cells



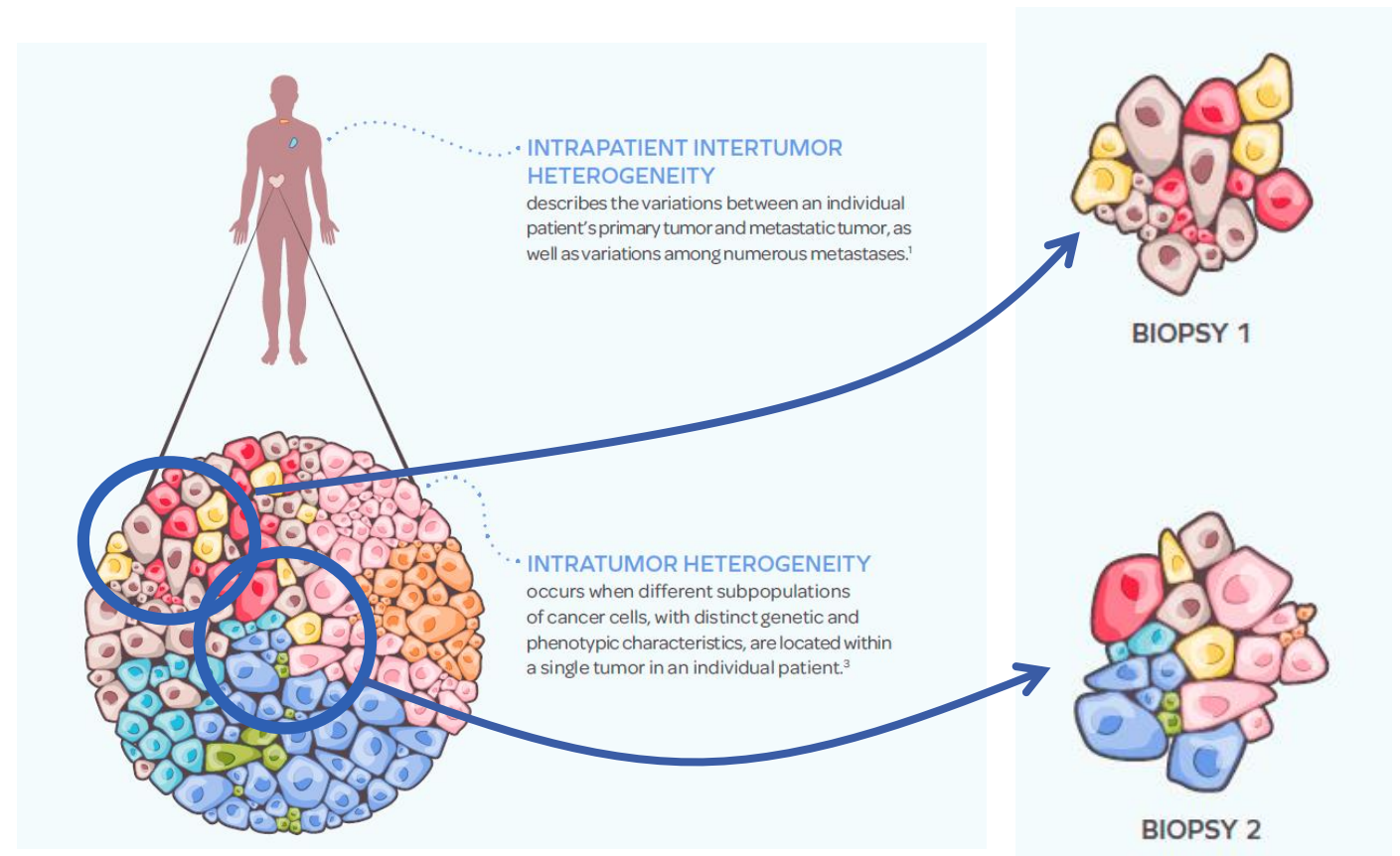
© 2014 Terese Winslow LLC
U.S. Govt. has certain rights

Why is SNV calling in cancer samples more complicated ?

Intra-tumour heterogeneity in which multiple tumour cell populations (subclones) exist

Multiple subclonal populations that are constantly evolving.

→ variants can be present in only one subclonal population



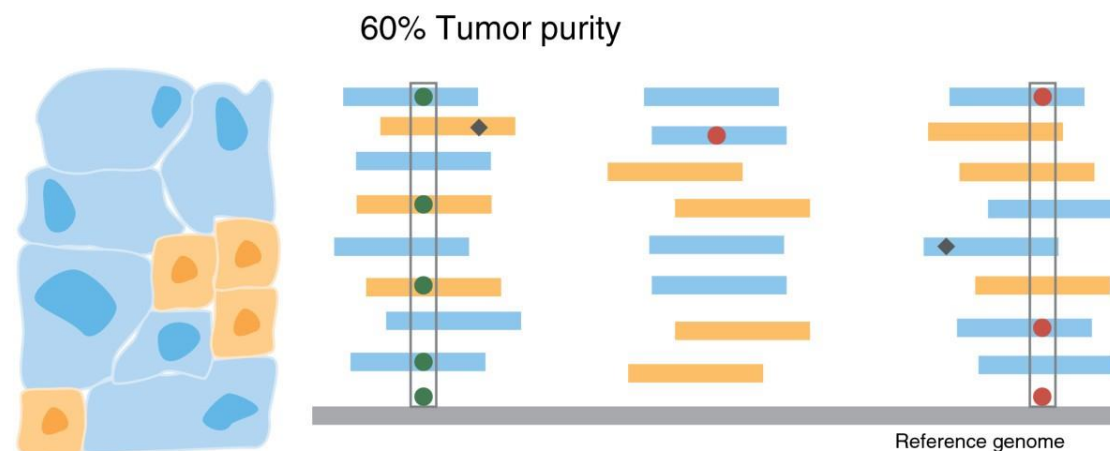
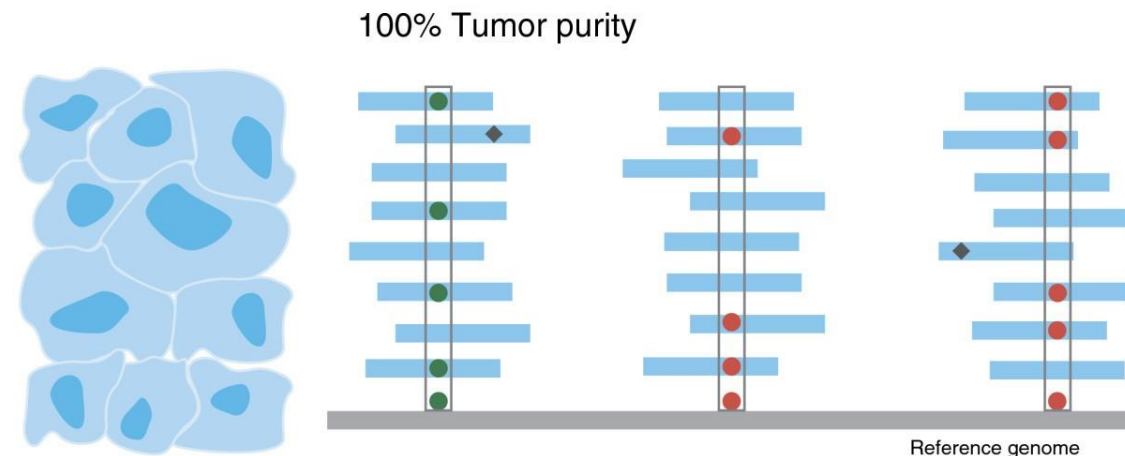
Why is SNV calling in cancer samples more complicated ?

Normal contamination

Normal cells can “contaminate” the tumor biopsy.

$$\textit{Tumor purity} = \frac{\textit{tumor cells}}{(\textit{normal} + \textit{tumor cells})}$$

Tumor sample



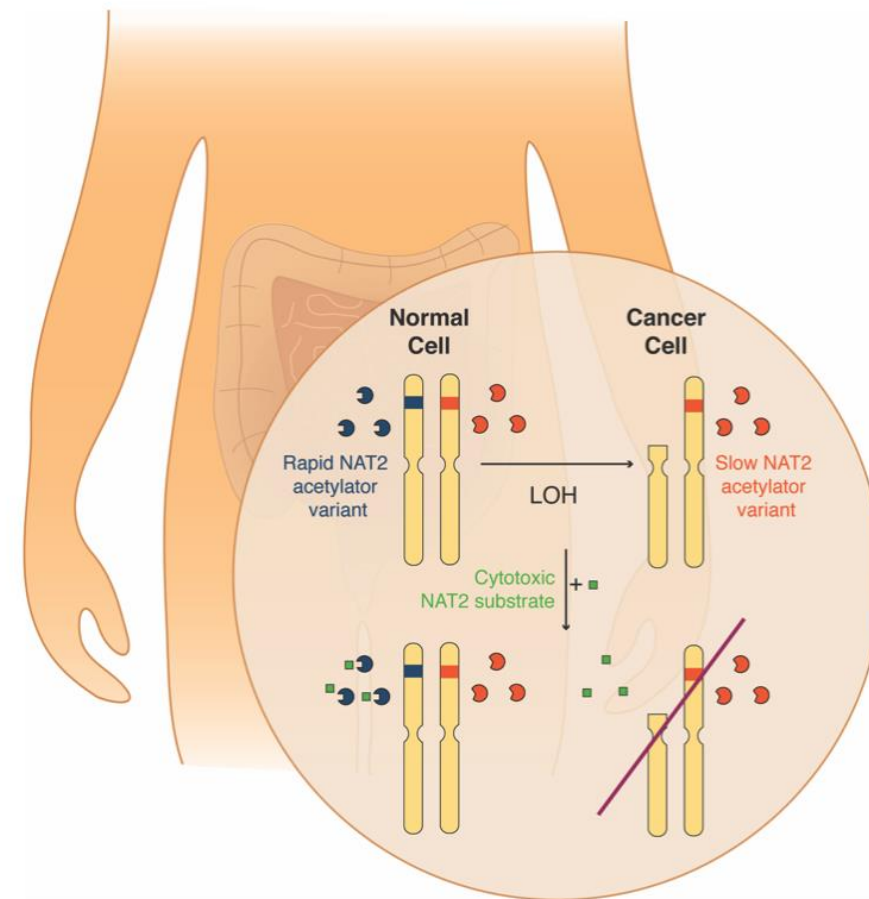
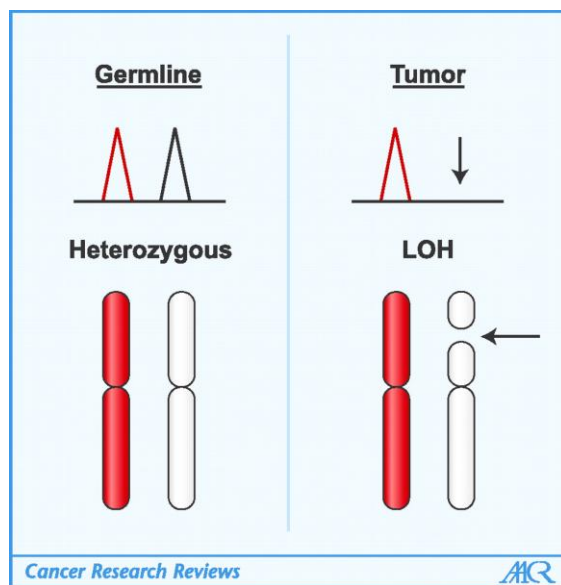
Key:

- Read
- Sequencing error
- Heterozygous germline SNV
- Heterozygous somatic SNV

Why is SNV calling in cancer samples more complicated?

Unbalanced structural variations (deletions, duplications, etc.)

→ need to detect LOH events

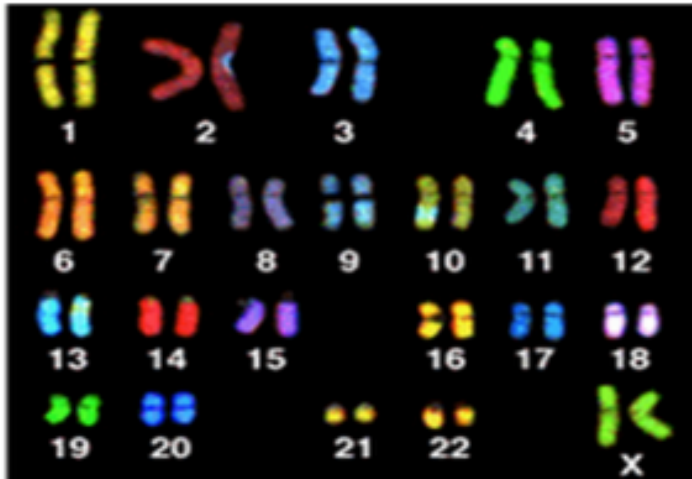


Why is SNV calling in cancer samples more complicated?

Aneuploidy → need for variant calling algorithm with no assumption on ploidy

Somatic alterations can be dramatic

Normal Cell



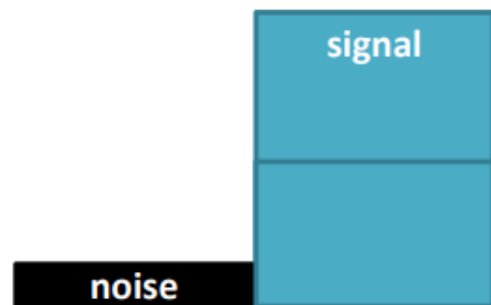
Cancer Cell Line HCC1954



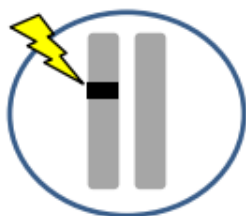
Why is SNV calling in cancer samples more complicated?

Amount of signal may be comparable to noise

Expectation for germline variants

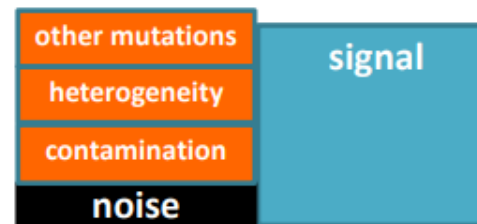


+ AF expected to follow ploidy

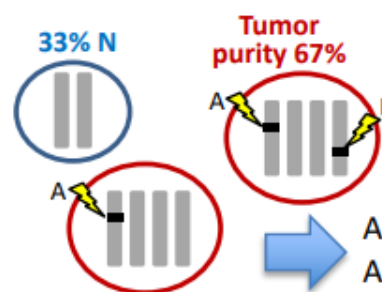


About 50% of reads will support the alternate allele

Expectation for somatic variants



+ no reliance on ploidy for AF

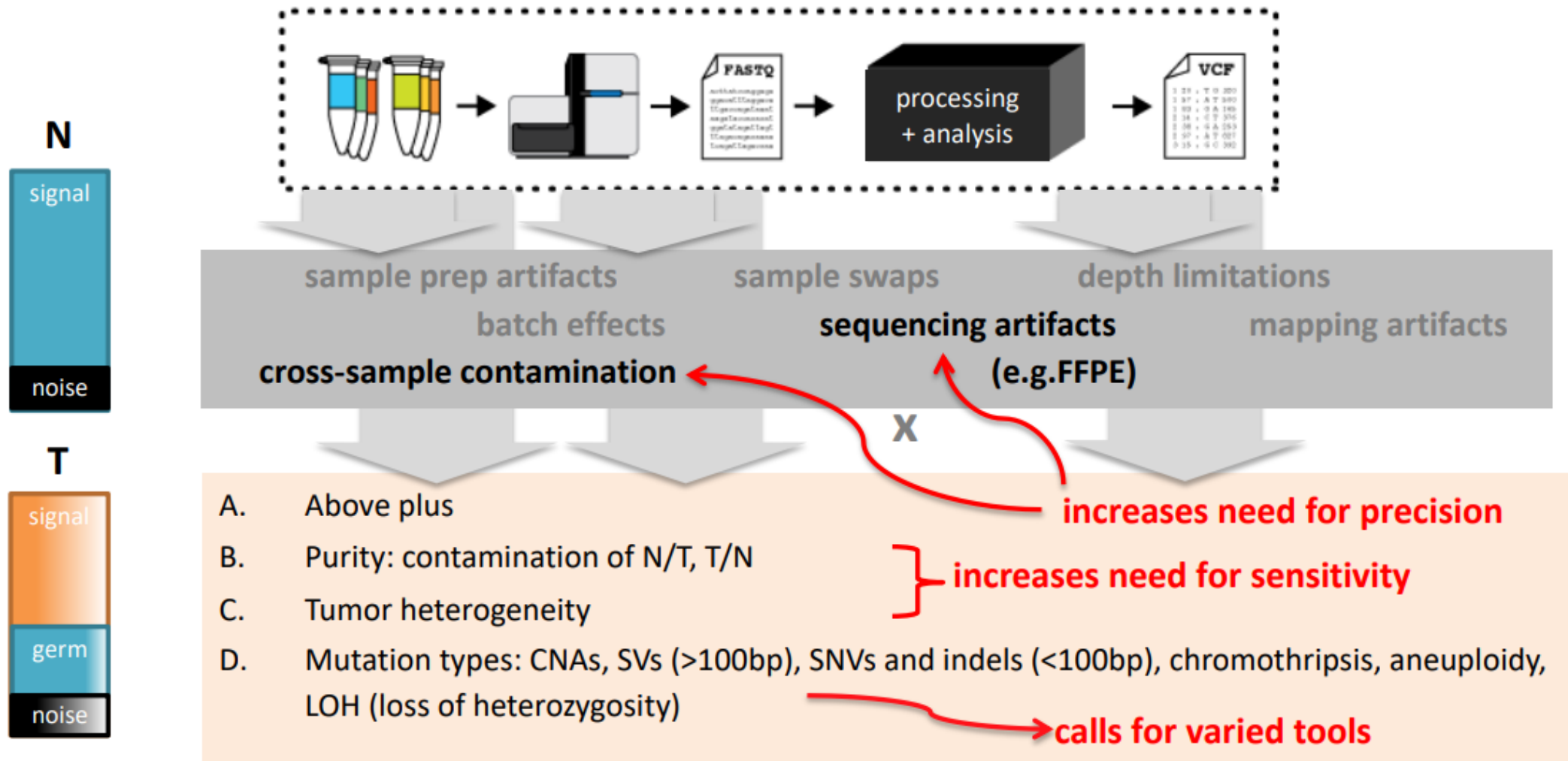


Example tumor sample:

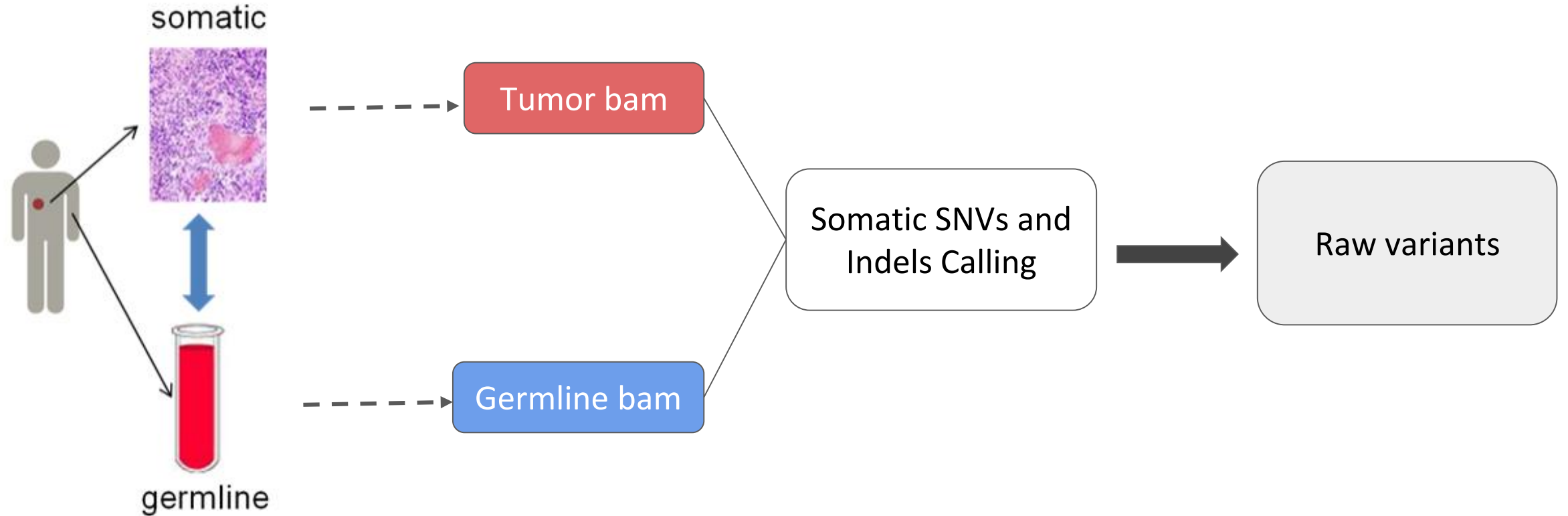
- Clonal SNV (A)
- Clonal copy number duplication
- Subclone of 50% of cancer cells (heterogeneity)
- Subclonal SNV (B)

About 20% of reads will support A
About 10% of reads will support B

Cancer-specific challenges confound analyses



Tumor - Normal pair analysis workflow

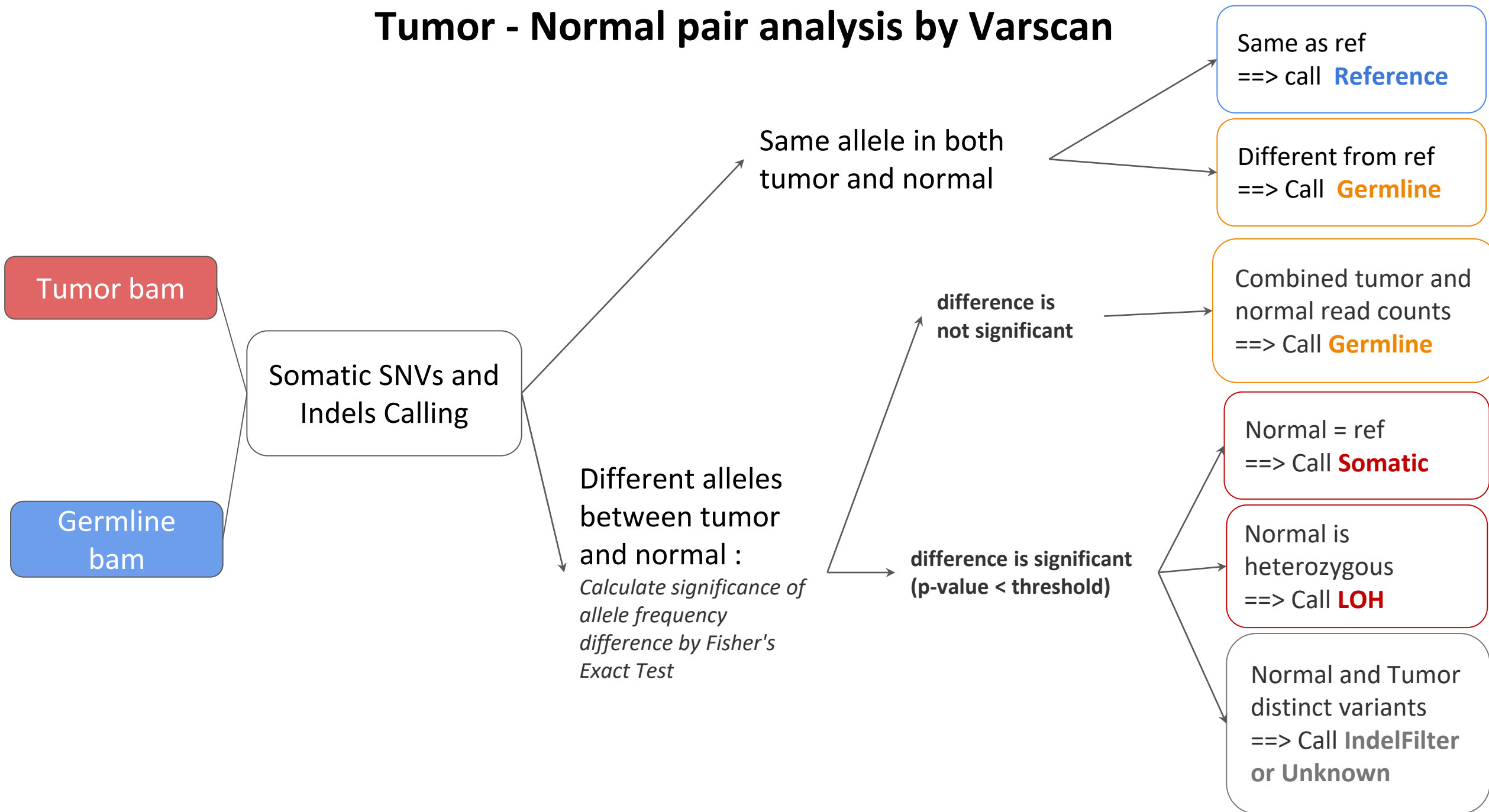


Logic of the Tumor-Normal workflow

Comparison to matched normal -> subtraction of germline background



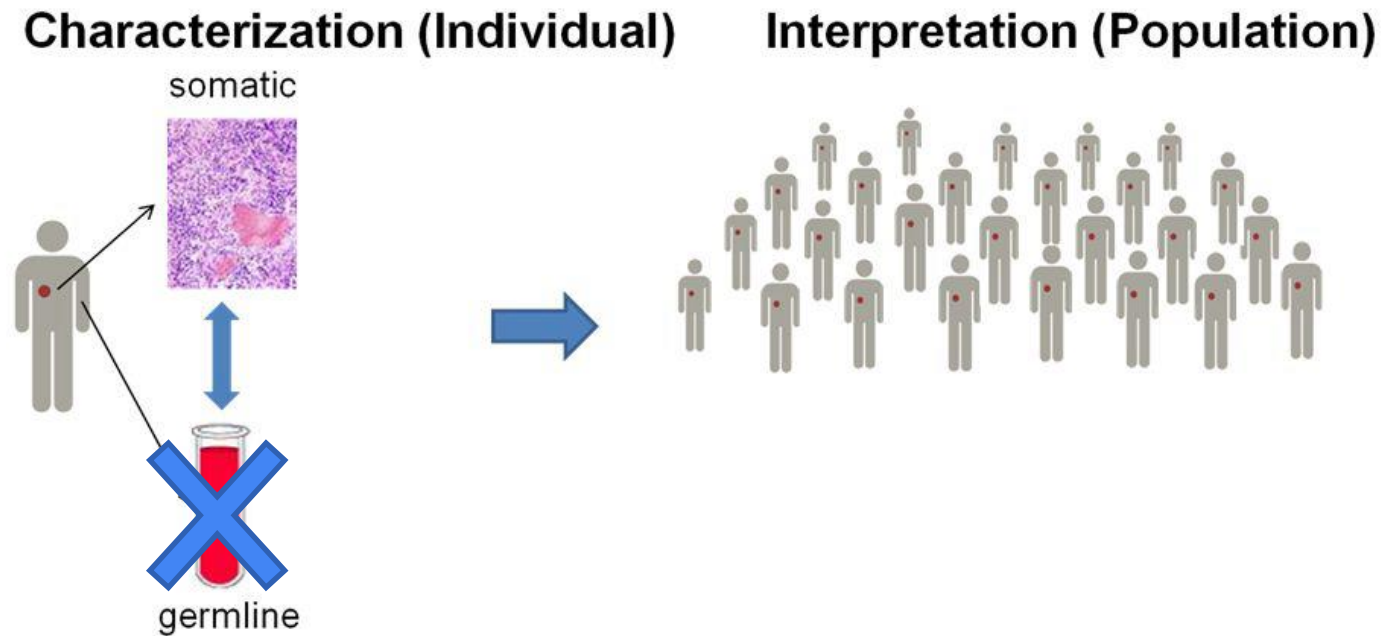
Tumor - Normal pair analysis by Varscan



Tumor-only analysis

If matched normal sample NOT available → Tumor only analysis

Pool of normals (PON) : used to eliminate common germline variation



Somatic variant filtering

Somatic variant callers output specific informations on variant :

→ Somatic likelihood like score as phred-scale somatic p-value

→ Status :
- Germline,
- Somatic,
- LOH

⇒ first metrics to filter variants on

Somatic variant annotation

Somatic Annotation Databases

Databases of variant-disease and gene-disease associations

- Cancer HotSpots
 - Single residue and in-frame indel mutation hotspots identified in 24,592 tumor samples.
- COSMIC
 - COSMIC (Catalogue of Somatic Mutations in Cancer) is a data resource that is designed to store and display somatic mutation information and related details and contains information relating to human cancers.
 - Data in COSMIC is curated from known Cancer Genes Literature and Systematic Screens.
- CIViC
 - CIViC (Clinical Interpretation of Variant in Cancer) is a an open access, open source, community-driven web resource for Clinical Interpretation of Variants in Cancer. The goal is to enable precision medicine by providing an educational forum for dissemination of knowledge and active discussion of the clinical significance of cancer genome alterations.



Somatic Annotation Databases

Databases of variant-disease and gene-disease associations

- Cancer Genome Interpreter (CGI)

Cancer Genome Interpreter (CGI) is designed to support the identification of tumor alterations that drive the disease and detect those that may be therapeutically actionable. CGI relies on existing knowledge collected from several resources and on computational methods that annotate the alterations in a tumor according to distinct levels of evidence.

It contains : a Cancer Biomarkers database, a Catalog of Validated Oncogenic Mutations and a Catalog of Validated Oncogenic Mutations

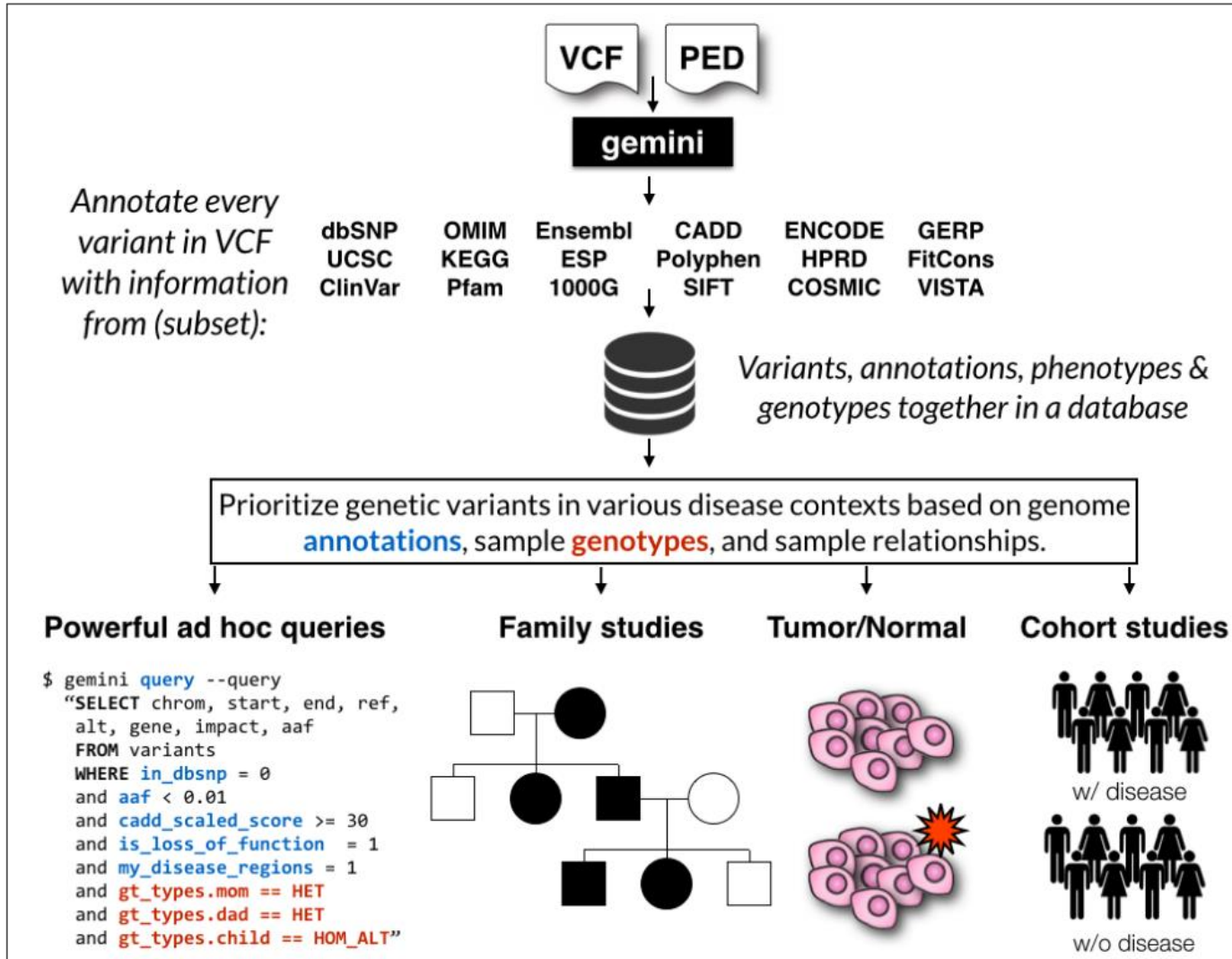
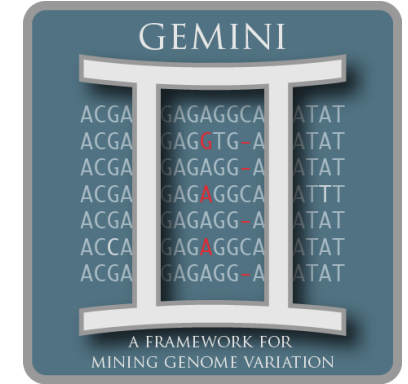


- The Cancer Genome Atlas (TCGA)

The Cancer Genome Atlas (TCGA), a landmark cancer genomics program, molecularly characterized over 20,000 primary cancer and matched normal samples spanning 33 cancer types. TCGA generated over 2.5 petabytes of genomic, epigenomic, transcriptomic, and proteomic data.



GEMINI presentation



Documentation :

<http://gemini.readthedocs.io>

Tutorials :

<https://speakerdeck.com/arq5x/>

GEMINI: a flexible framework for exploring genome variation

Overview

GEMINI (GEnome MINing) is a flexible framework for exploring genetic variation in the context of the wealth of genome annotations available for the human genome. By placing genetic variants, sample phenotypes and genotypes, as well as genome annotations into an integrated database framework, GEMINI provides a simple, flexible, and powerful system for exploring genetic variation for disease and population genetics.

Using the GEMINI framework begins by loading a VCF file (and an optional PED file) into a database. Each variant is automatically annotated by comparing it to several genome annotations from source such as ENCODE tracks, UCSC tracks, OMIM, dbSNP, KEGG, and HPRD. All of this information is stored in portable SQLite database that allows one to explore and interpret both coding and non-coding variation using "off-the-shelf" tools or an enhanced SQL engine.

Please also see the original manuscript.

Note

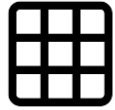
1. GEMINI solely supports human genetic variation mapped to build 37 (aka hg19) of the human genome.
2. GEMINI is very strict about adherence to VCF format 4.1.
3. For best performance, load and query GEMINI databases on the fastest hard drive to which you have access.

Tutorials

In addition to the documentation, please review the following tutorials if you are new to GEMINI. We recommend that you follow these tutorials in order, as they introduce concepts that build upon one another.

- Introduction to GEMINI, basic variant querying and data exploration. [html pdf](#)
- Identifying de novo mutations underlying Mendelian disease [html pdf](#)
- Identifying autosomal recessive variants underlying Mendelian disease [html pdf](#)
- Identifying autosomal dominant variants underlying Mendelian disease [html pdf](#)
- Other GEMINI tools [html pdf](#)

GEMINI database overview



The variants table

Gene information

gene	Genotype information		
transcript	gt		
is_exonic	gt_types		
is_coding	gt_phases		
is_lof	gt_depths		
is_splicing	gt_ref_depths		
exon	gt_alt_depths		
codon_change	gt_alt_freqs		
aa_change	gt_qual		
aa_length			
biotype			
impact			
impact_so			
impact_severity			
polyphen_pred			
polyphen_score			
sift_pred			
sift_score			
pfam_domain			

Variant and PopGen info

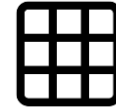
gt_types	type	Core VCF fields
gt_phases	sub_type	column_name
		chrom
		start
		end
gt_depths	call_rate	vcf_id
	num_hom_ref	variant_id
gt_ref_depths	num_het	anno_id
	num_hom_alt	ref
	num_unknown	alt
gt_alt_depths	aaf	qual
	hwe	filter
gt_alt_freqs	inbreeding_coeff	FLOAT The inbreeding coefficient
	pi	FLOAT The nucleotide diversity

Population information

in_dbsnp
rs_ids
in_hm2
in_hm3
in_esp
in_1kg
aaf_esp_ea
aaf_esp_aa
aaf_esp_all
aaf_1kg_amr
aaf_1kg_eas
aaf_1kg_sas
aaf_1kg_afr
aaf_1kg_eur
aaf_1kg_all
in_exac
aaf_exac_all
aaf_adj_exac_all
aaf_adj_exac_afr
aaf_adj_exac_amr
aaf_adj_exac_eas

Disease phenotype info (from ClinVar).

in_omim	BOOL	0 : Absence of the variant in OMIM databases 1 : Presence of the variant in OMIM databases
clinvar_causal_allele	STRING	The allele(s) that are associated or causal for the disease
clinvar_sig	STRING	The clinical significance scores for each of the variant according to ClinVar: <i>unknown, untested, non-pathogenic, probable-non-pathogenic, probable-pathogenic, pathogenic, drug-response, histocompatibility</i>
clinvar_disease_name	STRING	The name of the disease to which the variant is associated
clinvar_dbsource	STRING	Variant Clinical Channel IDs
clinvar_dbsource_id	STRING	The record id in the above database
clinvar_origin	STRING	The type of variant. Any of: <i>unknown, germline, somatic, inherited, paternal, maternal, de-novo, biparental, uniparental, not-tested, tested-inconclusive, other</i>
clinvar_dsdb	STRING	Variant disease database name
clinvar_dsdbid	STRING	Variant disease database ID
clinvar_disease_acc	STRING	Variant Accession and Versions
clinvar_in_locus_spec_db	BOOL	Submitted from a locus-specific database?
clinvar_on_diag_assay	BOOL	Variation is interrogated in a clinical diagnostic assay?
clinvar_gene_phenotype	STRING	' ' delimited list of phenotypes associated with this gene (includes any variant in the same gene in clinvar not just the current variant).



The variant_impacts table

variant_id	INTEGER	PRIMARY_KEY (Foreign key to <i>variants</i> table)
anno_id	INTEGER	PRIMARY_KEY (Based on variant transcripts)
gene	STRING	The gene affected by the variant.
transcript	STRING	The transcript affected by the variant.
is_exonic	BOOL	Does the variant affect an exon for this transcript?
is_coding	BOOL	Does the variant fall in a coding region (excludes 3' & 5' UTR's of exons)?
is_lof	BOOL	Based on the value of the impact col, is the variant LOF?
exon	STRING	Exon information for the variants that are exonic
codon_change	STRING	What is the codon change?
aa_change	STRING	What is the amino acid change?
aa_length	STRING	The length of CDS in terms of number of amino acids (SnpEff onLy)
biotype	STRING	The type of transcript (e.g., protein-coding, pseudogene, rRNA etc.) (SnpEff onLy)
impact	STRING	Impacts due to variation (ref.impact category)
impact_so	STRING	The sequence ontology term for the impact
impact_severity	STRING	Severity of the impact based on the impact column value (ref.impact category)
polyphen_pred	STRING	Impact of the SNP as given by PolyPhen (VEP onLy) benign, possibly_damaging, probably_damaging, unknown
polyphen_scores	FLOAT	Polyphen score reflecting severity (higher the impact, <i>higher</i> the score) (VEP onLy)
sift_pred	STRING	Impact of the SNP as given by SIFT (VEP onLy) neutral, deleterious
sift_scores	FLOAT	SIFT prob. scores reflecting severity (Higher the impact, <i>lower</i> the score) (VEP onLy)



Etc.

Tables/fields descriptions :

http://gemini.readthedocs.io/en/latest/content/database_schema.html

GEMINI usages

```
SELECT column-names  
FROM table-name  
WHERE condition  
ORDER BY sort-order
```



ad hoc data exploration

- Use SQL language to create queries and report data matching your requirements
- Can personalize your query to answer complex questions

```
gemini query -q "SELECT gene, chrom, clinvar_gene_phenotype FROM variants"
```

column_name	type
chrom	VARCHAR(20)
start	INTEGER
ref	TEXT
alt	TEXT
qual	FLOAT
filter	TEXT
in_omim	BOOLEAN
clinvar_sig	TEXT
clinvar_gene_phenotype	TEXT
gene	VARCHAR(60)

Table variants in Gemini database