

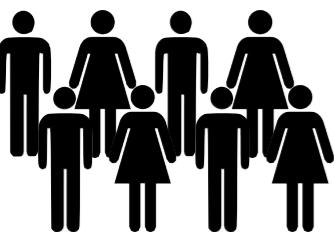


Variant interpretation and prioritization with GEMINI

Aaron Quinlan | University of Virginia | Oct. 3, 2013
Beyond the Genome, 2013 | San Francisco, CA

Typical genetics study designs

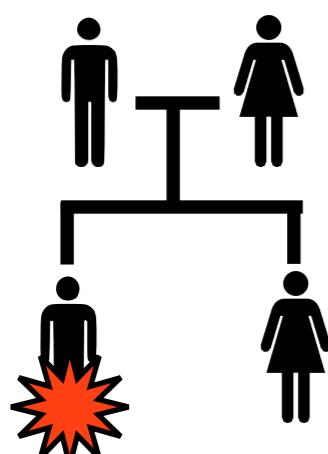
disease



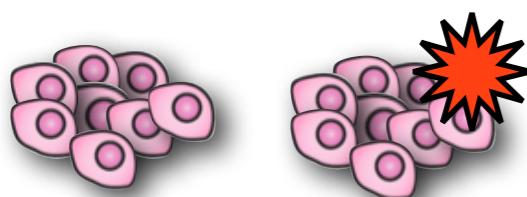
no disease



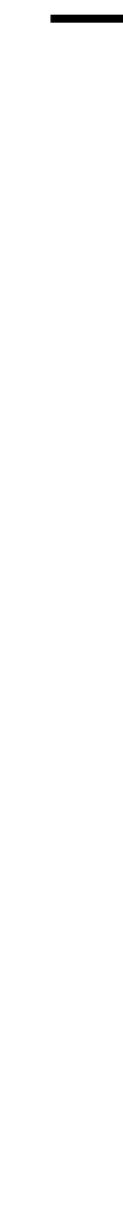
“case” v. “control”



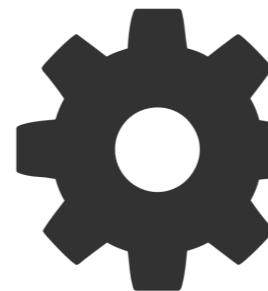
family studies



cancer genomics

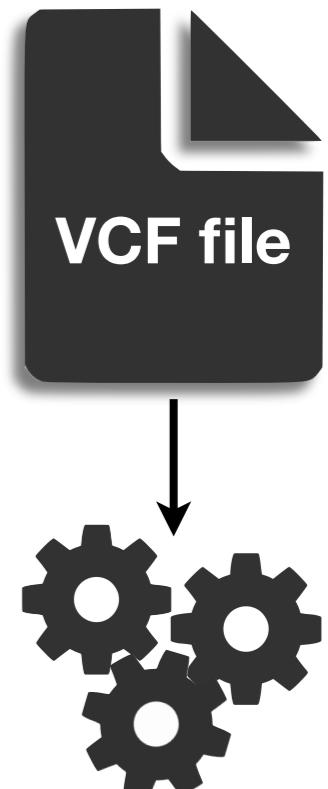
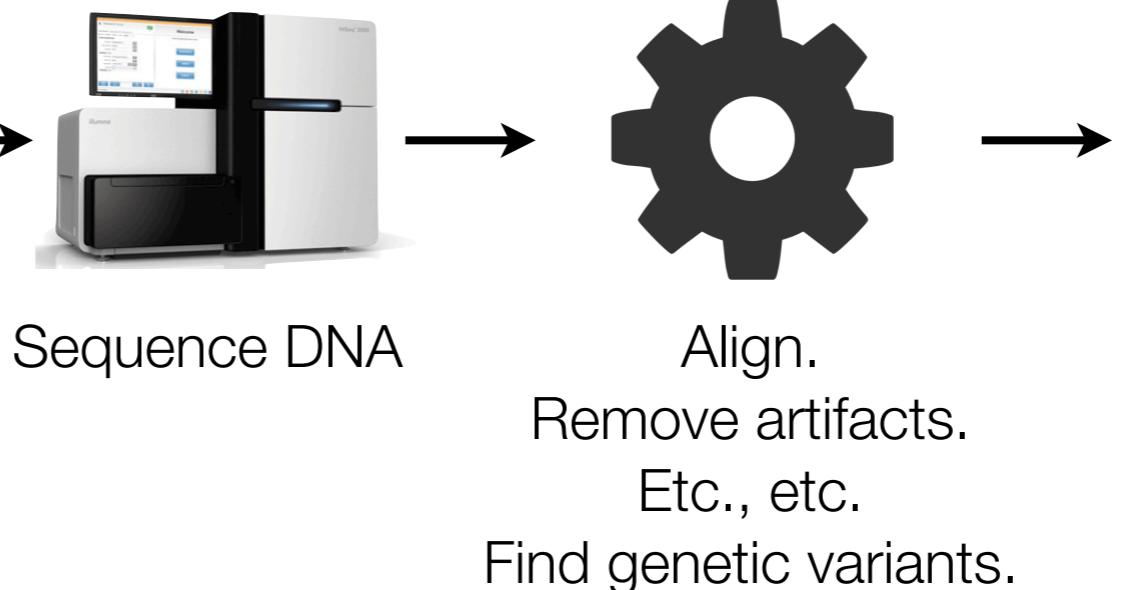
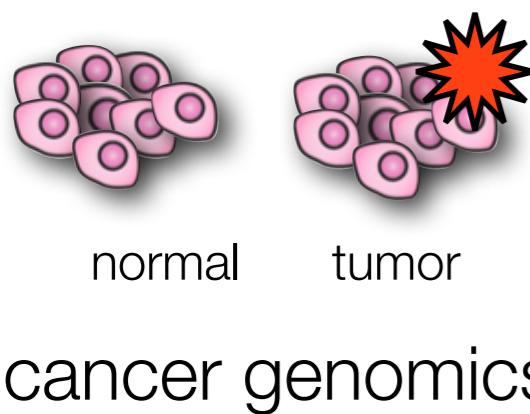
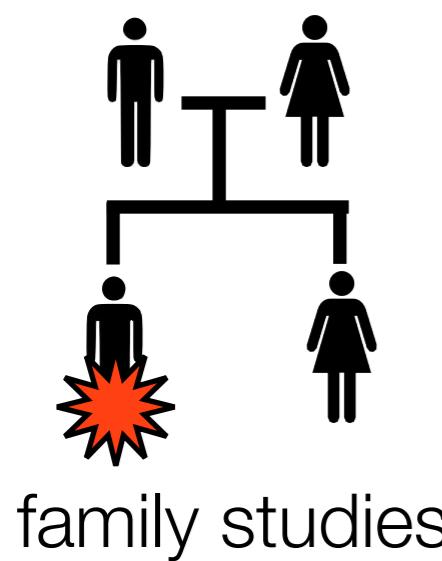
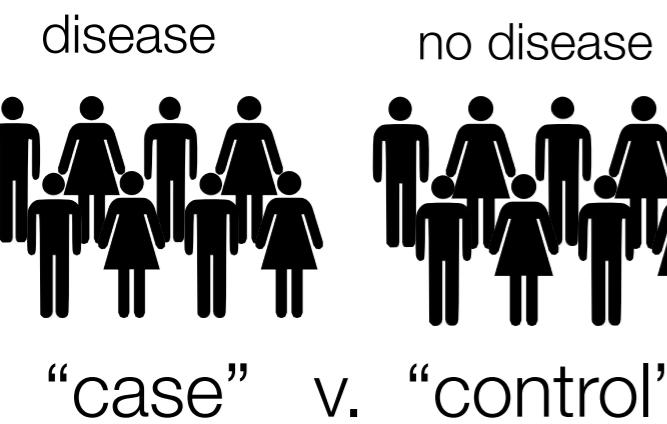


Sequence DNA



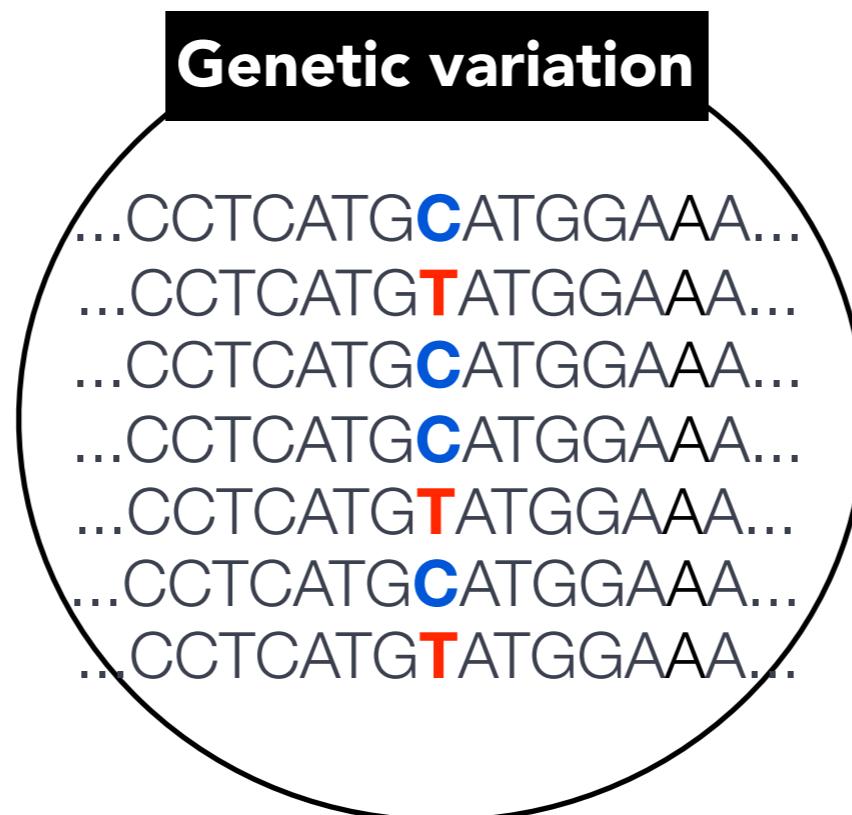
Align.
Remove artifacts.
Etc., etc.
Find genetic variants.

Typical genetics study designs

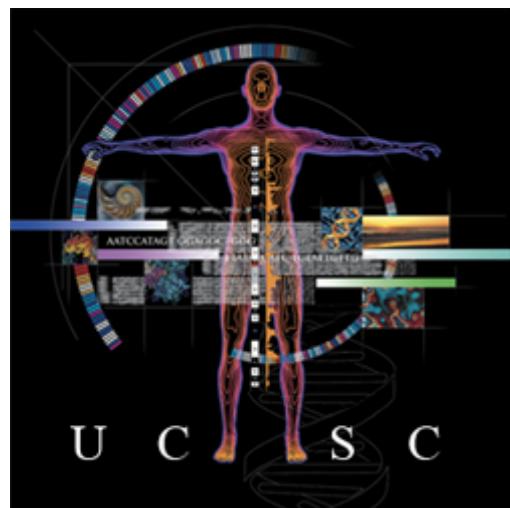


Interpret.
Prioritize.
Repeat.

Analytical challenges: data integration

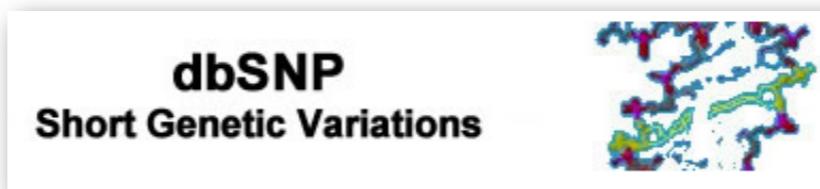


Analytical challenges: data integration

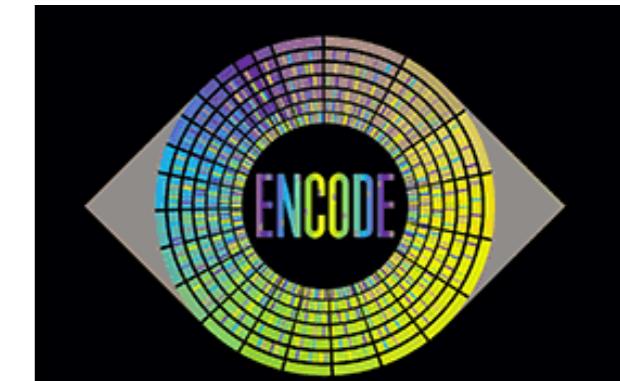


Conservation
Repeat elements
Genome Gaps
Cytobands
Gene annotations
"Mappability"
DeCIPHER
ISGA

Pfam



ClinVar



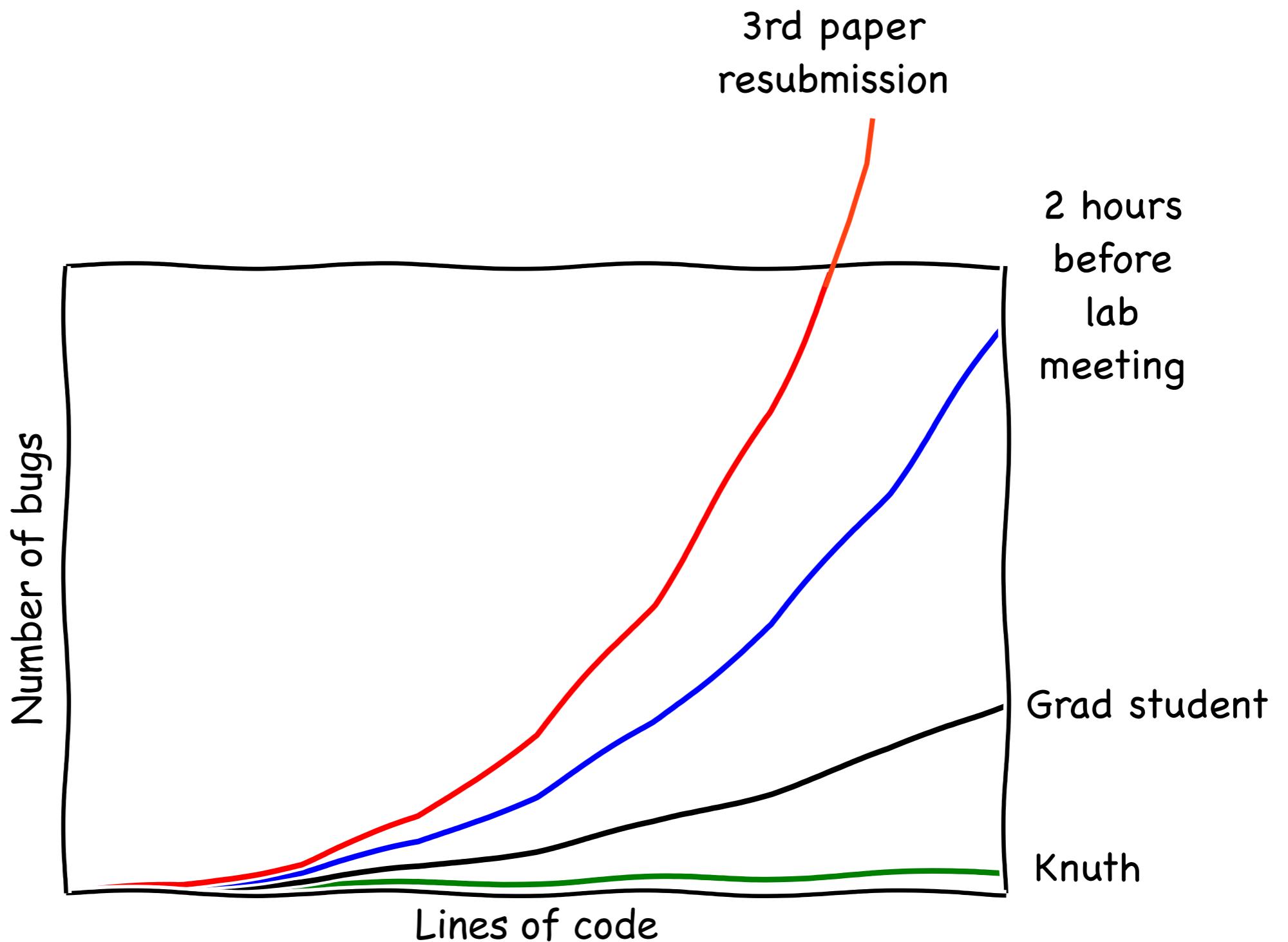
Genetic variation

...CCTCATGCATGGAAA...
...CCTCATGTATGGAAA...
...CCTCATGCATGGAAA...
...CCTCATGCATGGAAA...
...CCTCATGTATGGAAA...
...CCTCATGCATGGAAA...
...CCTCATGTATGGAAA...

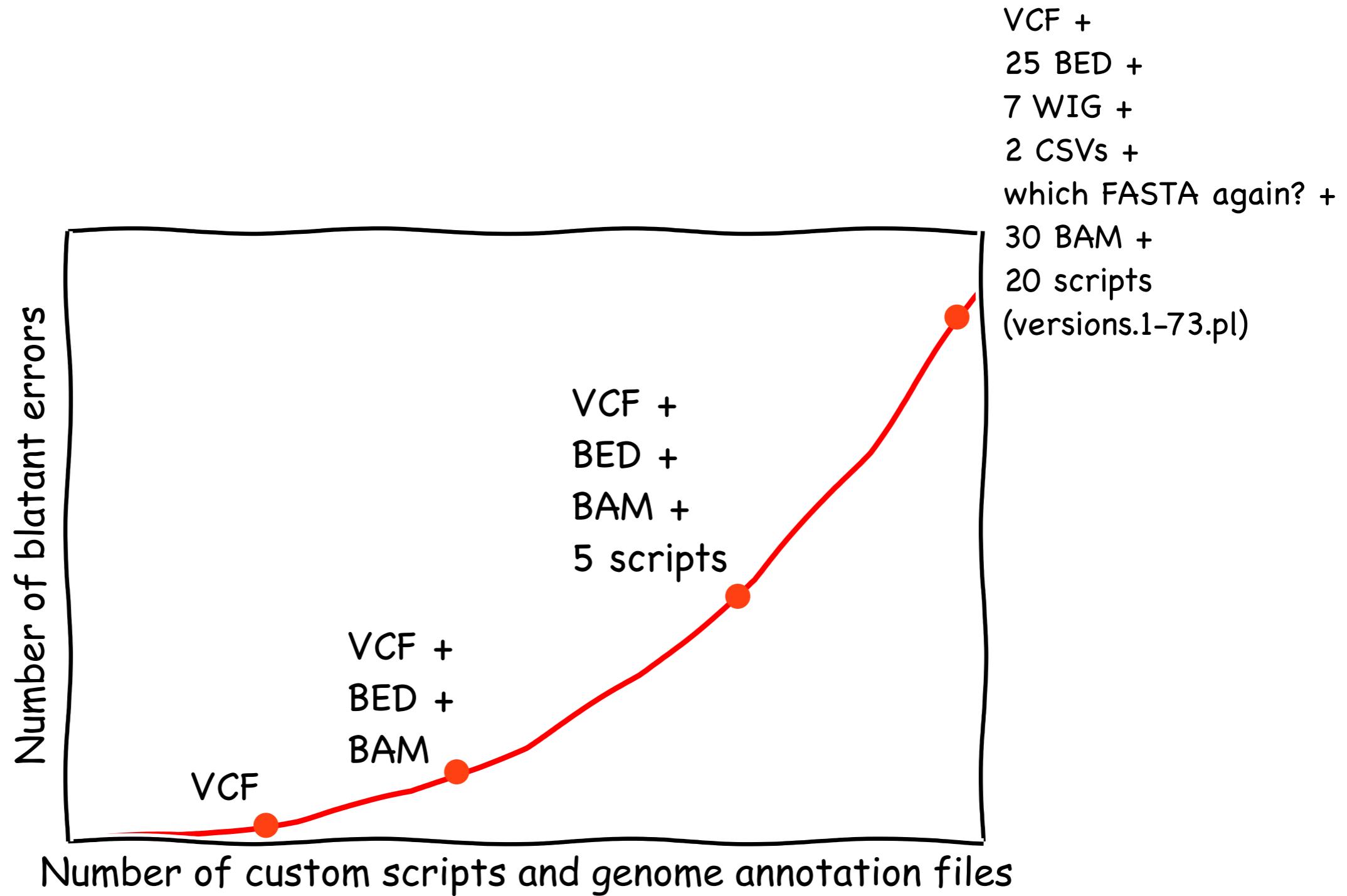
Chromatin marks
DNA methylation
RNA expression
TF binding

More code, more problems

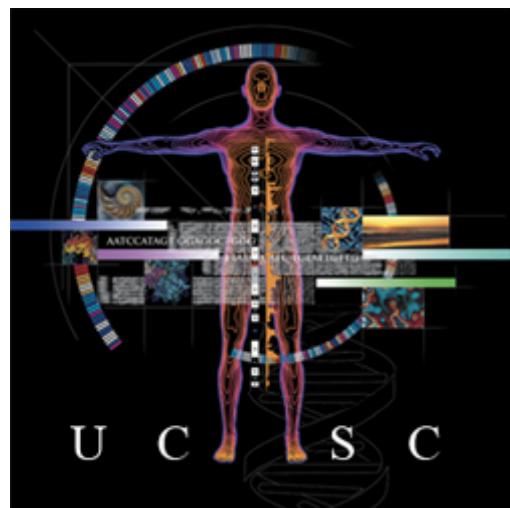
More code, more problems



A genomics corollary



Directly Integrate variants with annotations

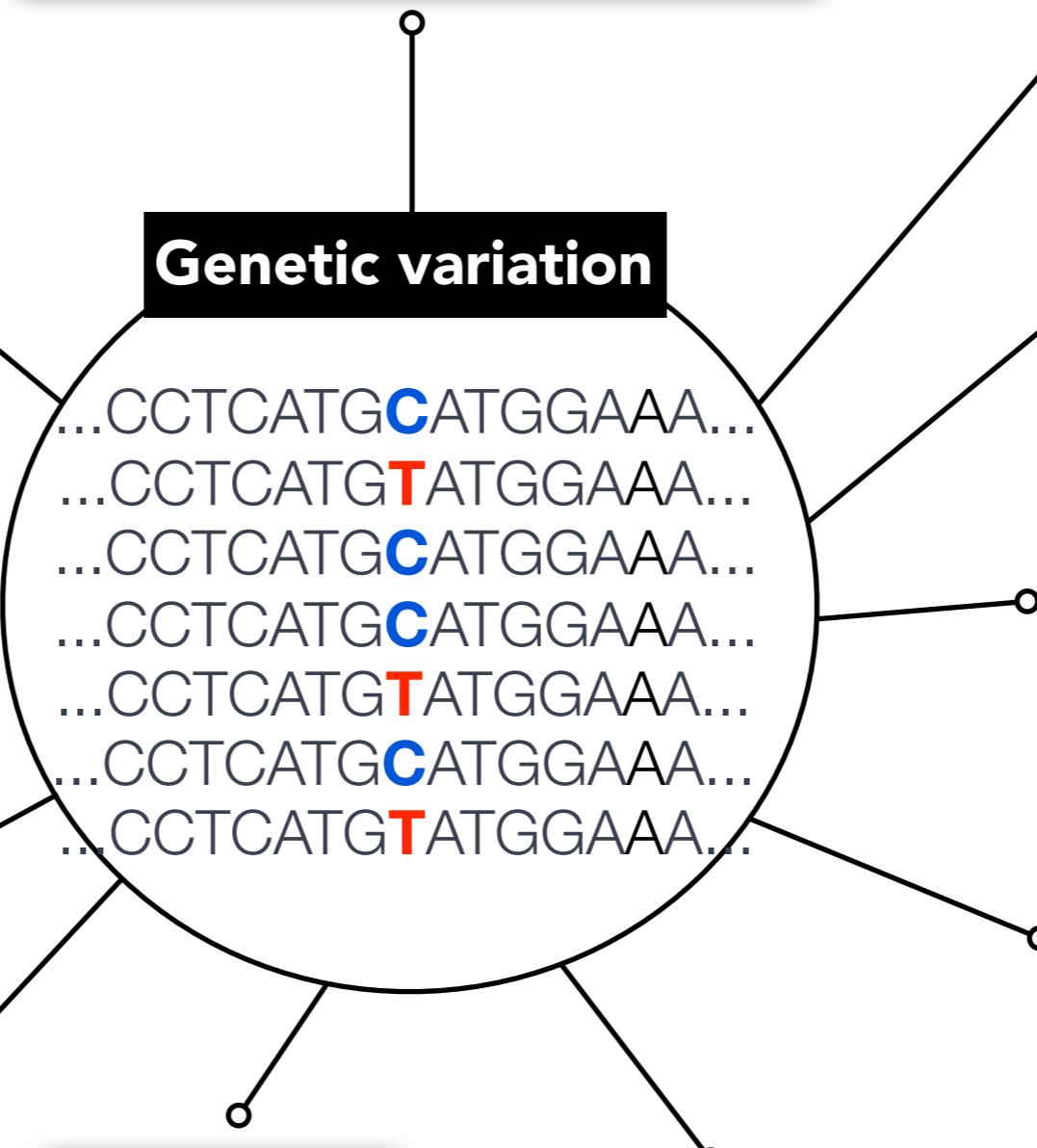
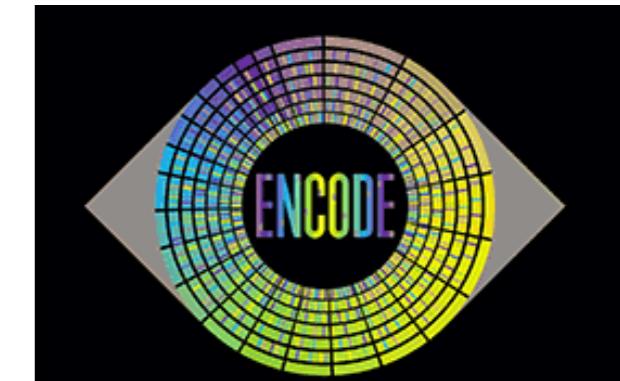


Conservation
Repeat elements
Genome Gaps
Cytobands
Gene annotations
"Mappability"
DeCIPHER
ISGA

Pfam



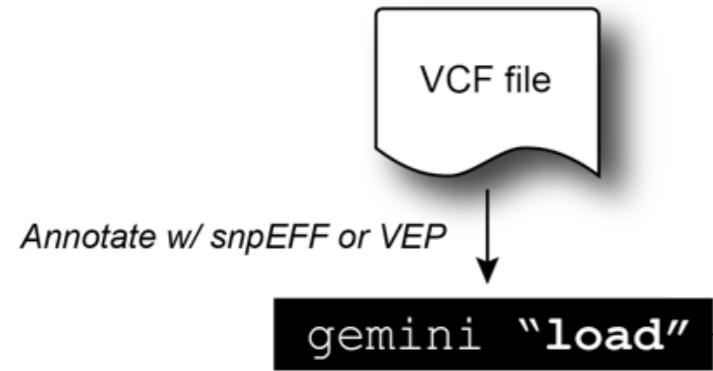
ClinVar



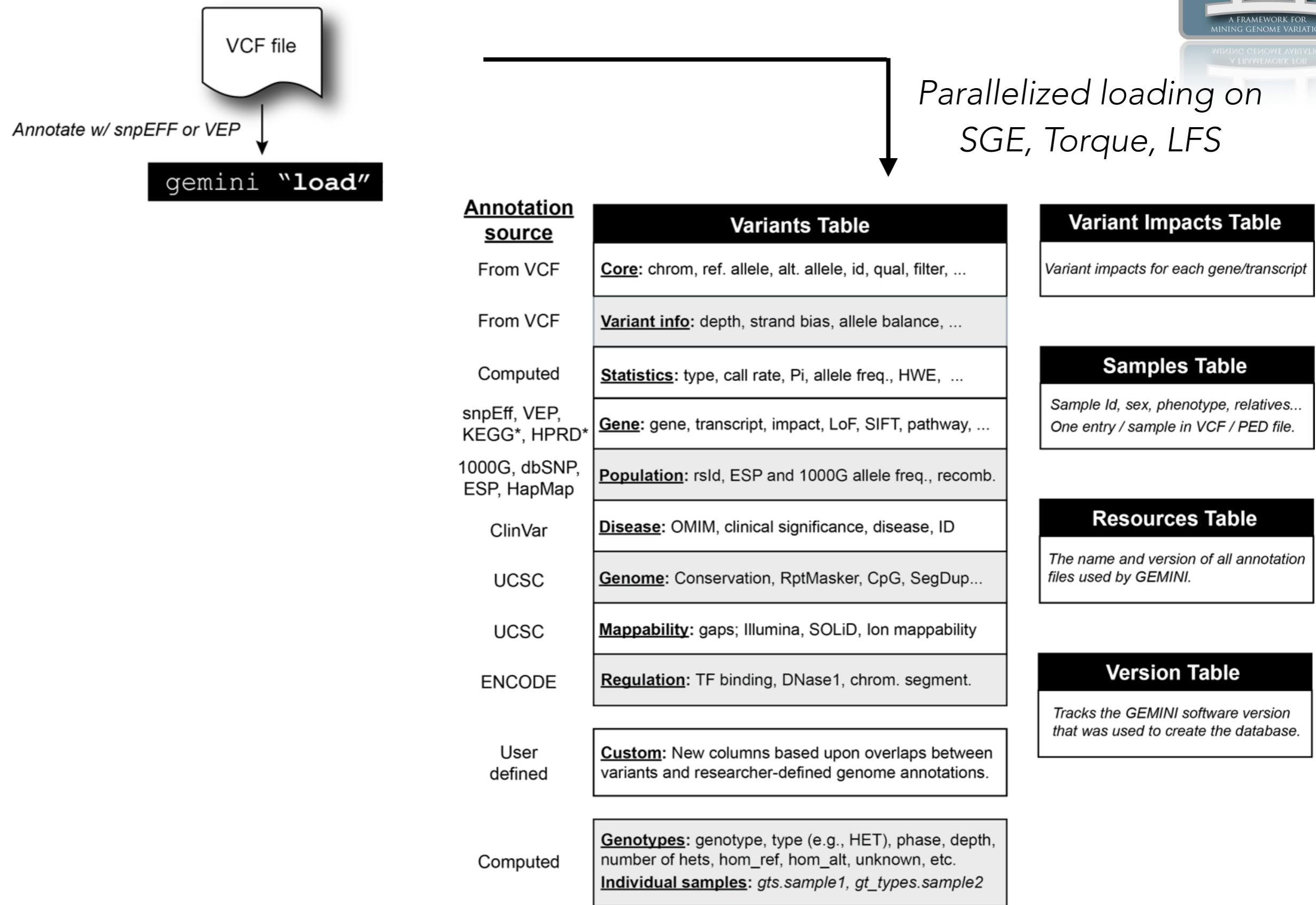
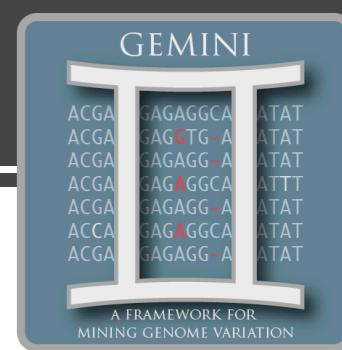
Chromatin marks
DNA methylation
RNA expression
TF binding



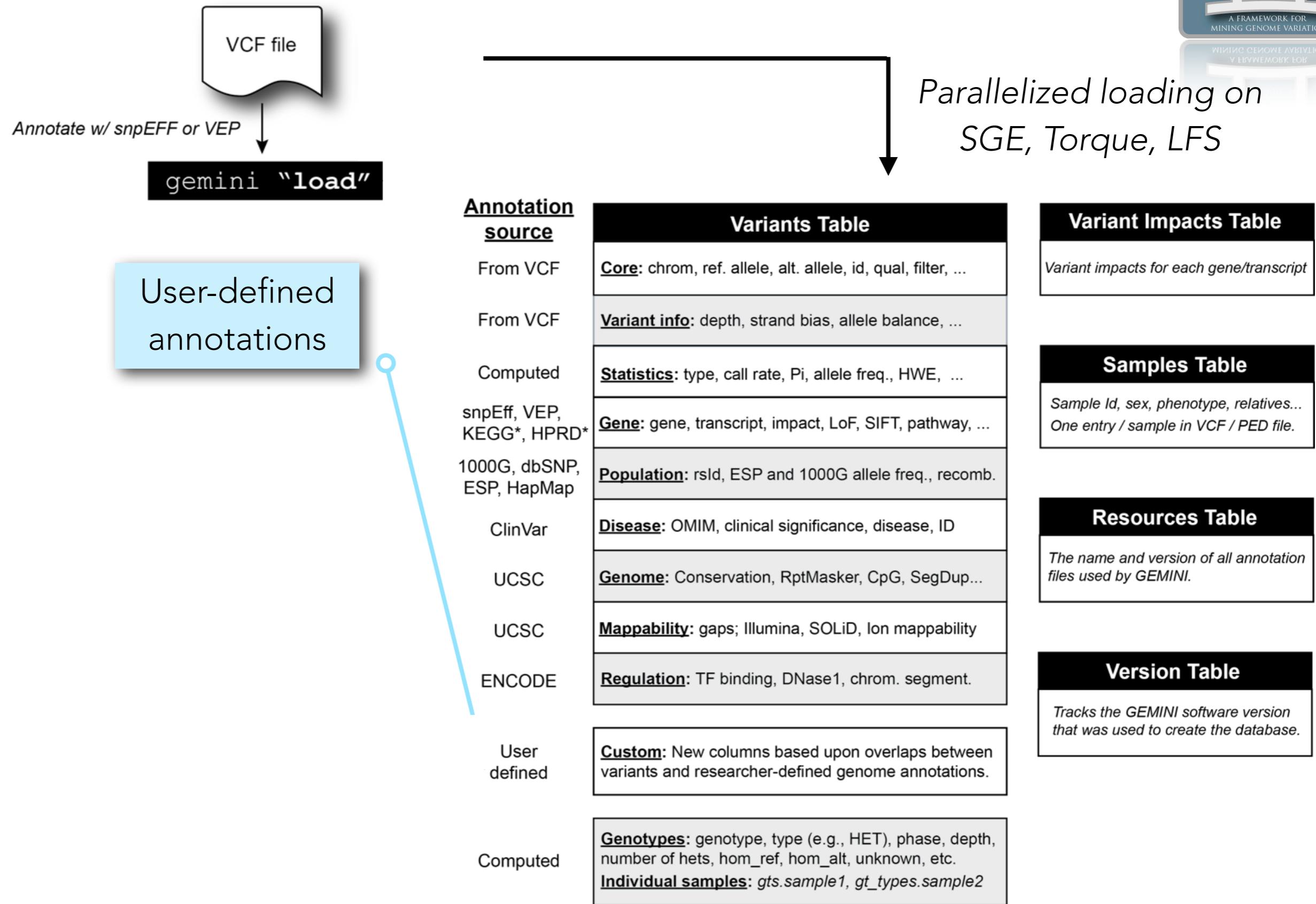
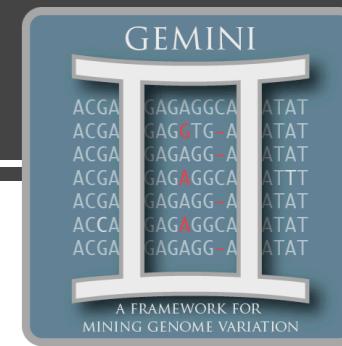
The GEMINI framework



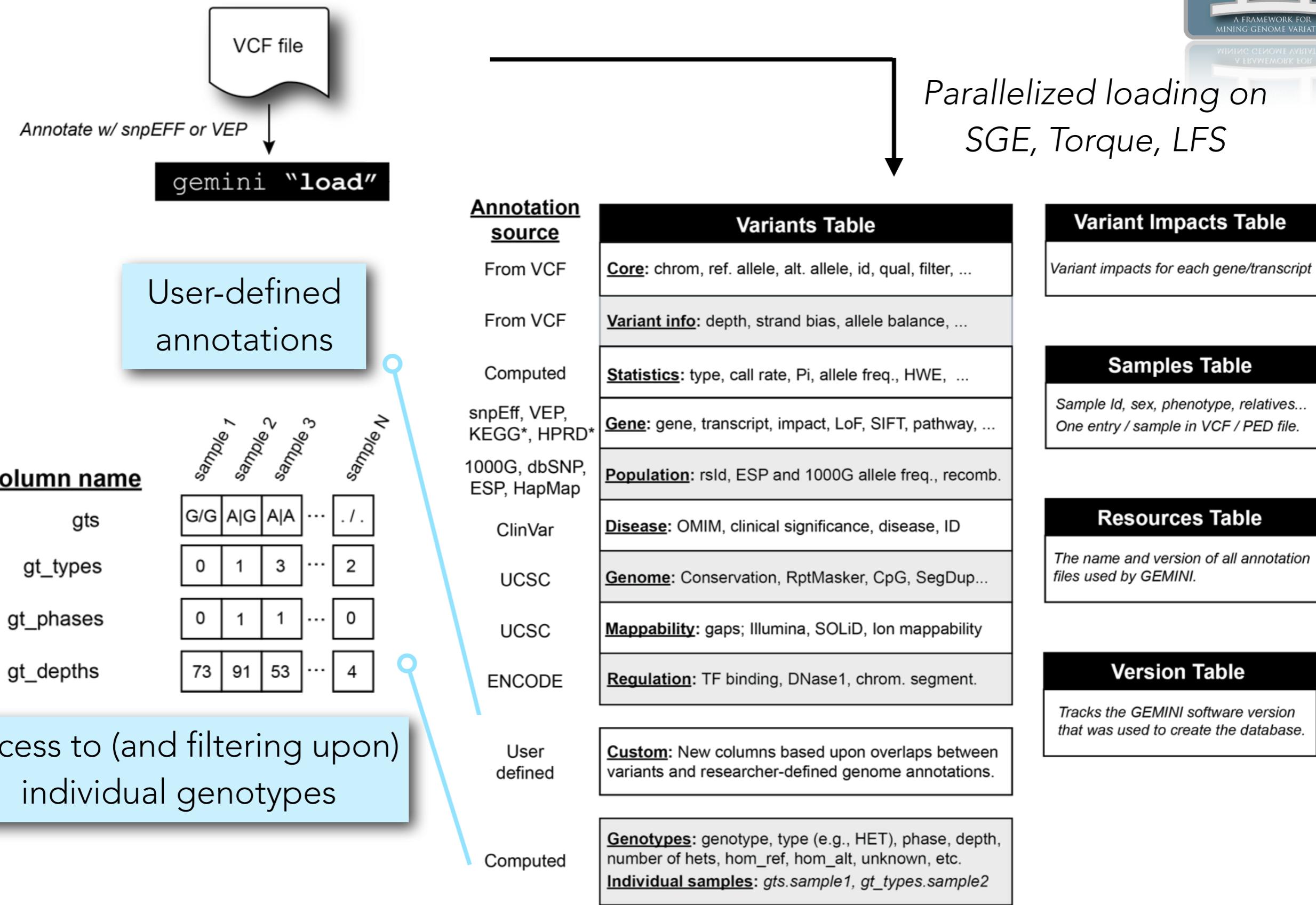
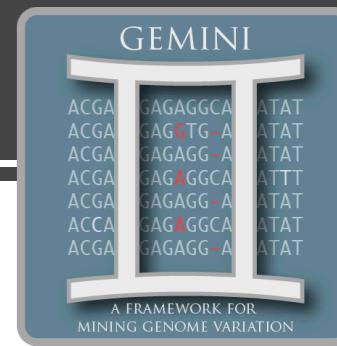
The GEMINI framework



The GEMINI framework



The GEMINI framework



We can exploit sample *relationships* (not in VCF)

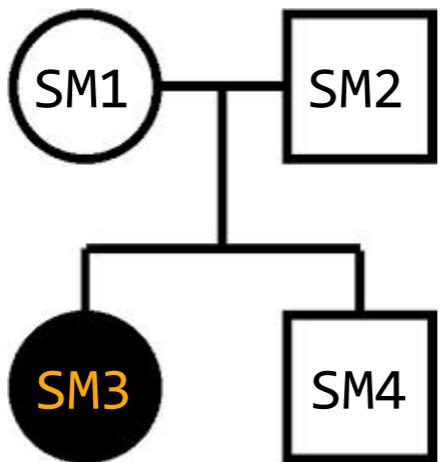


		SM1	SM2	SM3	SM4
CHR1	123	A	G	A/G	A/G
				G/G	A/A

We can exploit sample *relationships* (not in VCF)



			SM1	SM2	SM3	SM4	
CHR1	123	A	G	A/G	A/G	G/G	A/A

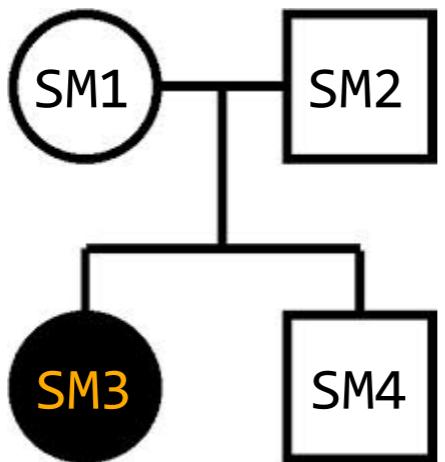


PED file

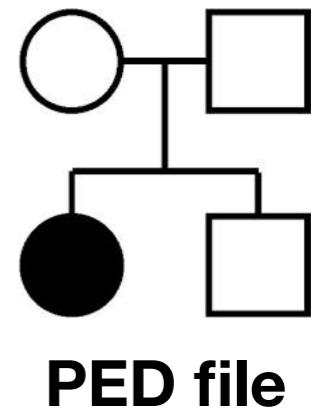
We can exploit sample relationships (not in VCF)



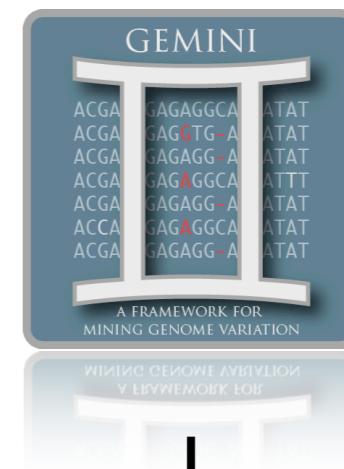
		SM1	SM2	SM3	SM4
CHR1	123	A	G	A/G	A/G



PED file



PED file

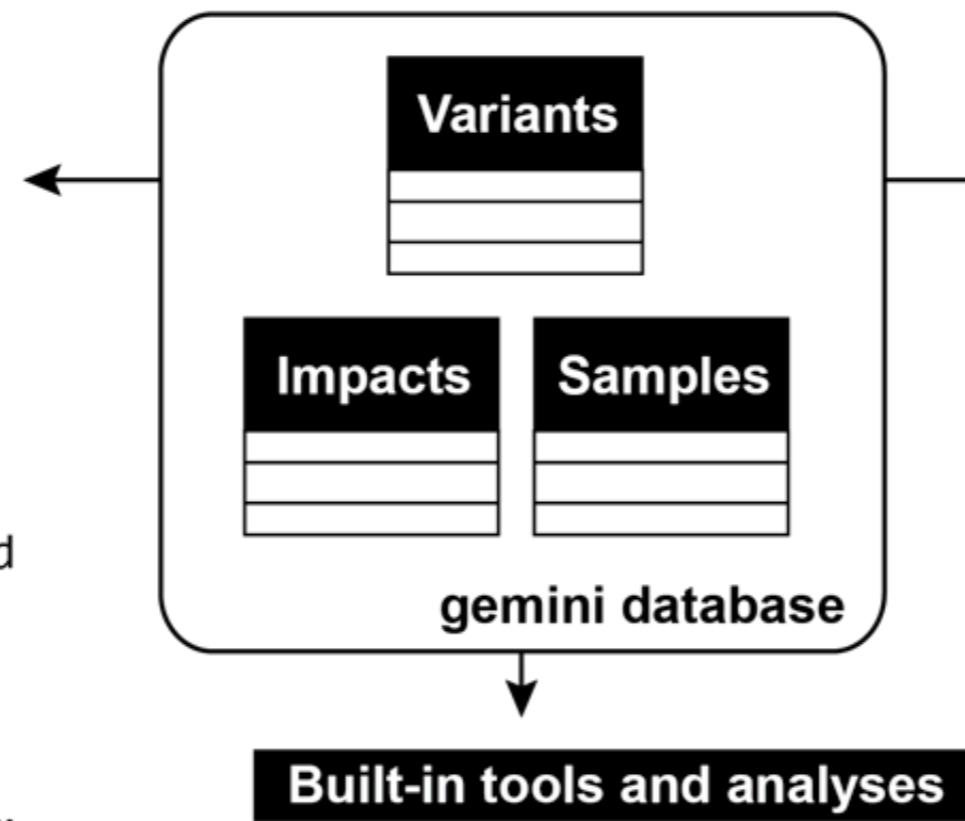


Disease models
Inheritance patterns
Transmission tests
Phasing

Mining variation with GEMINI

ad hoc data exploration

```
gemini query  
--query  
"select chrom, start, end,  
ref, alt, gene,  
impact, aaf, gts.kid  
from variants  
where in_dbsnp = 0  
and aaf < 0.01  
and is_lof = 1  
and my_disease_regions = 1"  
  
--gt-filter  
"gt_types.mom == HET  
and  
gt_types.dad == HET  
and  
gt_types.proband == HOM_ALT"
```



Framework for new tools

- Burden tests
- Population genetics
- Pedigree studies
- Haplotype analysis
- Custom tools and new methods

Built-in tools and analyses

Tool	Description
region	extract variants spec. genomic intervals or genes
stats	compute variant statistics (SFS, Ts/Tv, counts, etc.)
annotate	add new columns based on custom annotations
windower	compute variant statistics across genome “windows”
comp_hets	identify candidate compound heterozygotes
pathways	maps genes and variants to KEGG pathways
lof_sieve	prioritize candidate loss-of-function variants
interact	find protein interactions for genes/variants/samples
auto_rec	identify variants meeting an autosomal recessive model
auto_dom	identify variants meeting an autosomal dominant model
de_novo	identify candidate de novo mutations
browser	launch the interactive gemini web browser interface

Example analyses: *ad hoc* exploration

1. Find all rare (<1%), LoF variants

```
$ gemini query -q "select * from variants  
    where is_lof = 1  
    and aaf <= 0.01"  
my.gemini.db
```

Example analyses: *ad hoc* exploration

1. Find all rare (<1%), LoF variants

```
$ gemini query -q "select * from variants  
    where is_lof = 1  
    and aaf <= 0.01"  
my.gemini.db
```

2. Find all LoF variants and show me genotypes
for all *affected samples having blue eyes*

```
$ gemini query -q "select * from variants  
    where is_lof = 1"  
    --sample-filter "phenotype=1 and eye_color='blue'"  
my.gemini.db
```



Automatically built from custom sample
annotations in PED file

Example analyses: *ad hoc* exploration

3. Find variants with a known clinical phenotype

```
$ gemini query -q "select * from variants  
    where (in_omim = 1  
        or clinvar_disease_name is not NULL)"
```

Example analyses: *ad hoc* exploration

3. Find variants with a known clinical phenotype

```
$ gemini query -q "select * from variants
  where (in_omim = 1
    or clinvar_disease_name is not NULL)"
  --gt-filter "gt_types.proband == HET
  and
  gt_types.mommy == HOM_REF
  and
  gt_types.daddy == HOM_REF"
```

Example analyses: *ad hoc* exploration

3. Find variants with a known clinical phenotype

```
$ gemini query -q "select * from variants
  where (in_omim = 1
    or clinvar_disease_name is not NULL)"
  --gt-filter "gt_types.proband == HET
  and
  gt_types.mommy == HOM_REF
  and
  gt_types.daddy == HOM_REF"
  --min-kindreds 3
```

Example analyses: *ad hoc* exploration

3. Find variants with a known clinical phenotype

```
$ gemini query -q "select * from variants
  where (in_omim = 1
    or clinvar_disease_name is not NULL)"
  --gt-filter "gt_types.proband == HET
  and
  gt_types.mommy == HOM_REF
  and
  gt_types.daddy == HOM_REF"
  --min-kindreds 3
  --region chr1:149000000-153000000
```

Example analyses: *built-in analyses*

1. Find *de novo* mutations

```
$ gemini de_novo my.gemini.db
```

2. Find compound heterozygotes

```
$ gemini comp_hets my.gemini.db
```

3. Find variants meeting an auto. recessive pattern

```
$ gemini auto_recessive my.gemini.db
```

4. Find variants meeting an auto. dominant pattern

```
$ gemini auto_dominant my.gemini.db
```

Example: autosomal recessive candidates

gemini **auto_recessive**

--columns “chrom, start, end, ref,
alt, gene, impact”

request specific columns
from database

--filter “impact_severity != ‘LOW’”

restrict to higher impact
variants

--min-kindreds 2

found in >= 2 families

--depth 40

each sample must have
>= 40 reads

disease.gemini.db

fam	genotypes (F,M,C)	depths	gene	chrom	start	end	ref	alt	impact
1	T/C,T/C,C/C	69,48,70	WDR37	chr10	1142207	1142208	T	C	stop_loss
2	T/C,T/C,C/C	59,49,82	WDR37	chr10	1142207	1142208	T	C	stop_loss

Ongoing research.

Burden tests

		Genetic variants											
		Binary trait											
		Gene 1	Gene 2	Gene N	Gene 1	Gene 2	Gene N	Gene 1	Gene 2	Gene N	Gene 1	Gene 2	Gene N
Cases		1	1	1	1	0	1	1	0	1	1	0	1
Controls		1	1	1	1	0	0	1	1	1	1	0	1
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													
.													

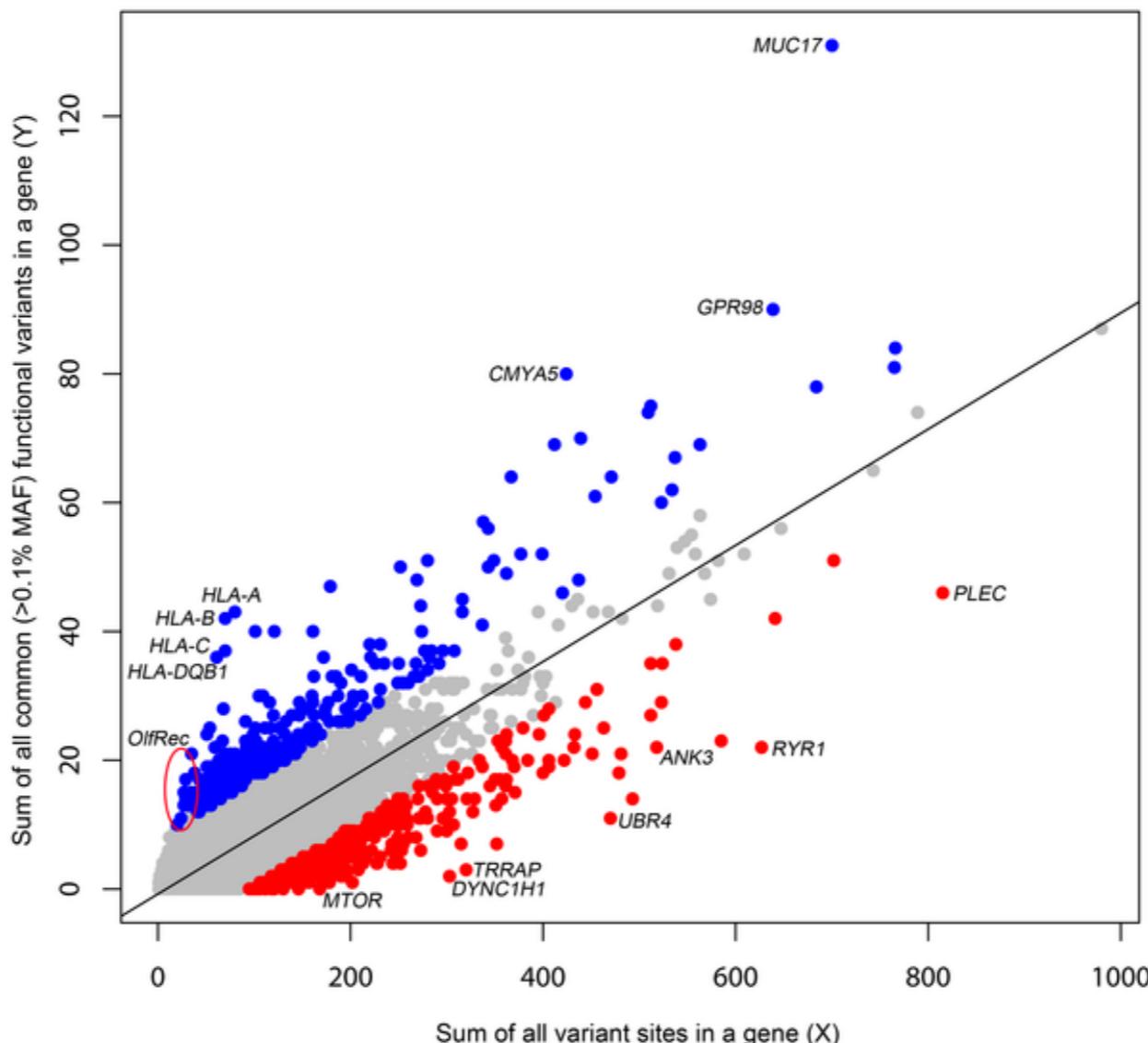
Gene prioritization

OPEN  ACCESS Freely available online

 PLOS GENETICS

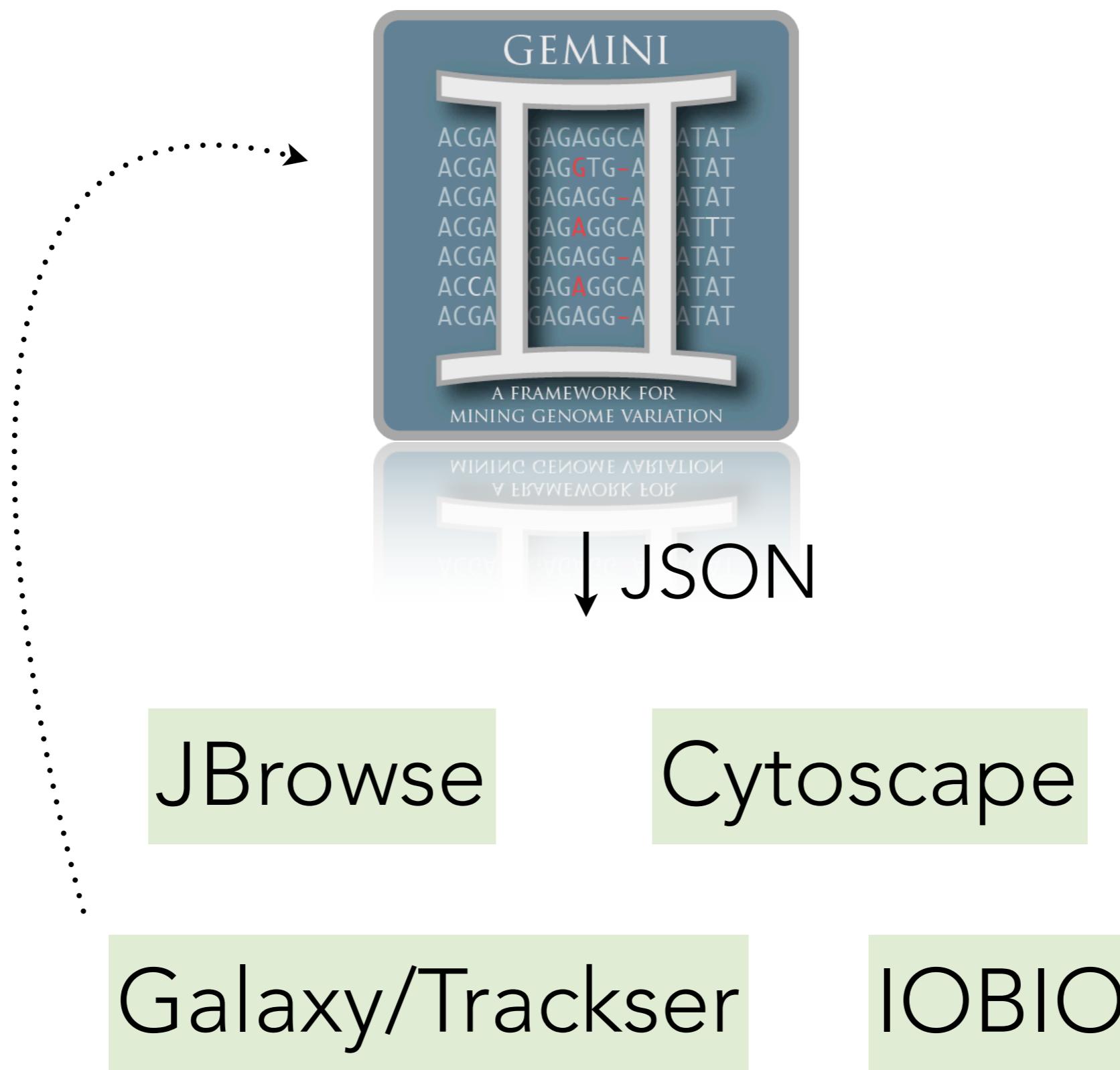
Genic Intolerance to Functional Variation and the Interpretation of Personal Genomes

Slavé Petrovski^{1,2*}, Quanli Wang¹, Erin L. Heinzen^{1,3}, Andrew S. Allen^{1,4}, David B. Goldstein^{1*}



An API for data viz. and tool dev.

query



Summary

- Flexible framework for mining genetic variation.
- Integrates important genome annotations.
- Query access to individual genotypes
- Extensible for new analyses and tool dev.
- Free. Open source. github.com/arg5x/gemini
- Well documented. gemini.readthedocs.org
- Extensible, portable, & reproducible

OPEN  ACCESS Freely available online

 PLOS COMPUTATIONAL
BIOLOGY

GEMINI: Integrative Exploration of Genetic Variation and Genome Annotations

Umadevi Paila¹, Brad A. Chapman², Rory Kirchner², Aaron R. Quinlan^{1*}

Acknowledgements



HARVARD | **SCHOOL OF PUBLIC HEALTH**
Powerful ideas for a healthier world



Uma Paila*

Postdoctoral Fellow

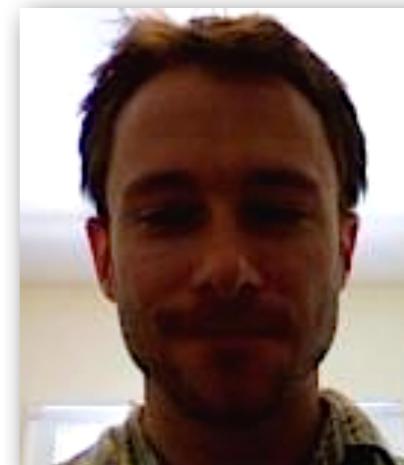
github.com/udp3f



Brad Chapman



Oliver Hofmann



Rory Kirchner