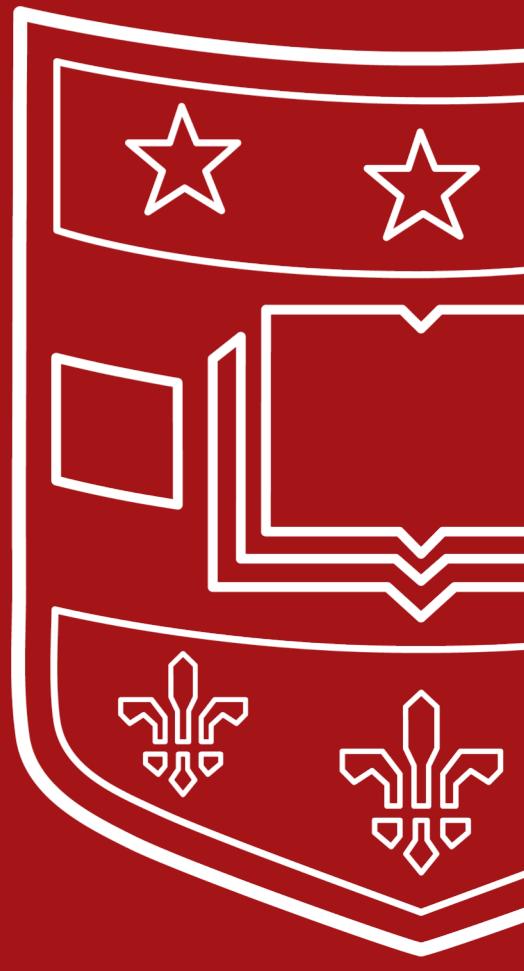


Clinical Interpretation Knowledgebases and Learning Bioinformatics

Alex H Wagner, PhD
December, 2019 – UTH BIO390

 Washington University in St. Louis





Part I – Clinical Interpretation Knowledgebases



Swiss Institute
of
Bioinformatics

Introduction to bioinformatics: Clinical Bioinformatics

Valérie Barbié, head of SIB Clinical Bioinformatics

Zürich, 26 November 2019



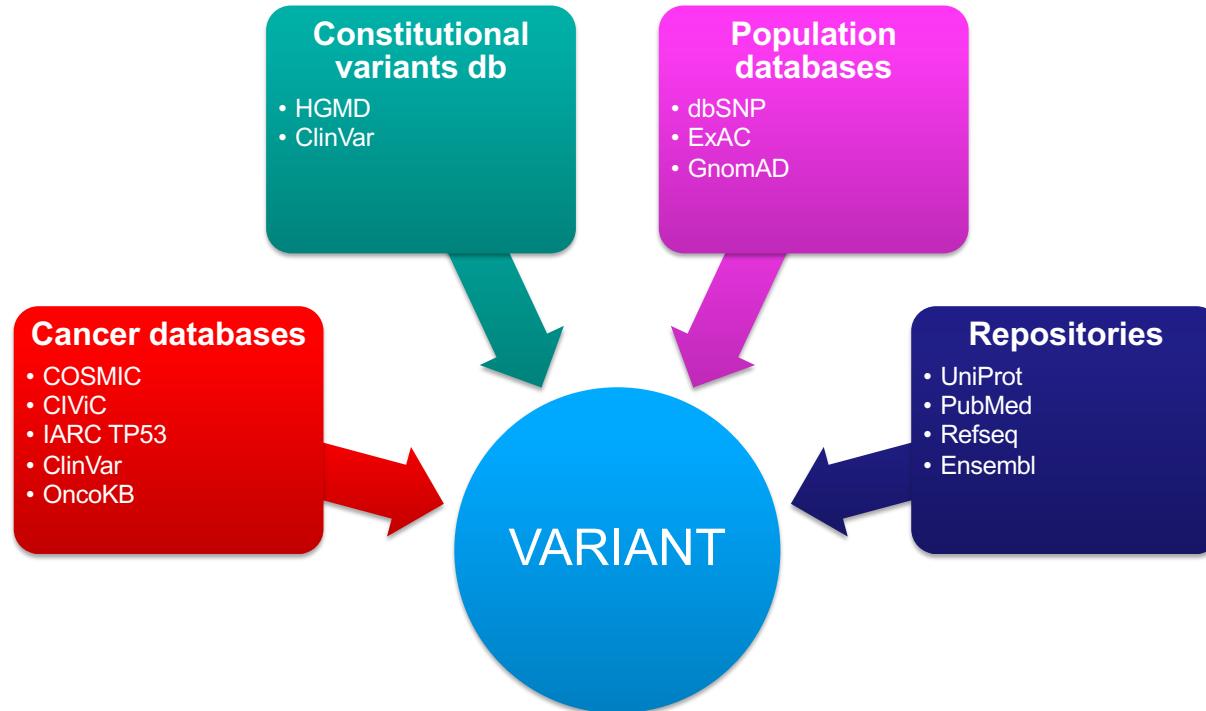
www.sib.swiss



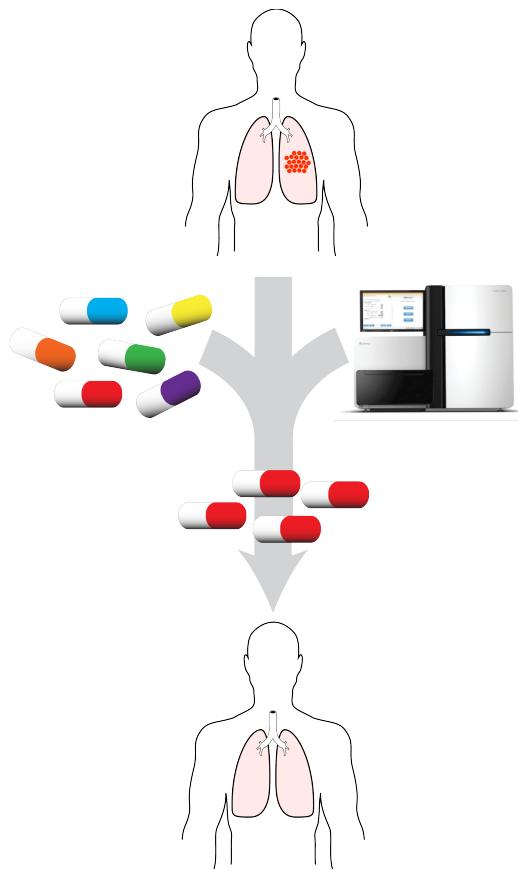


Annotating a variant: knowledgebases

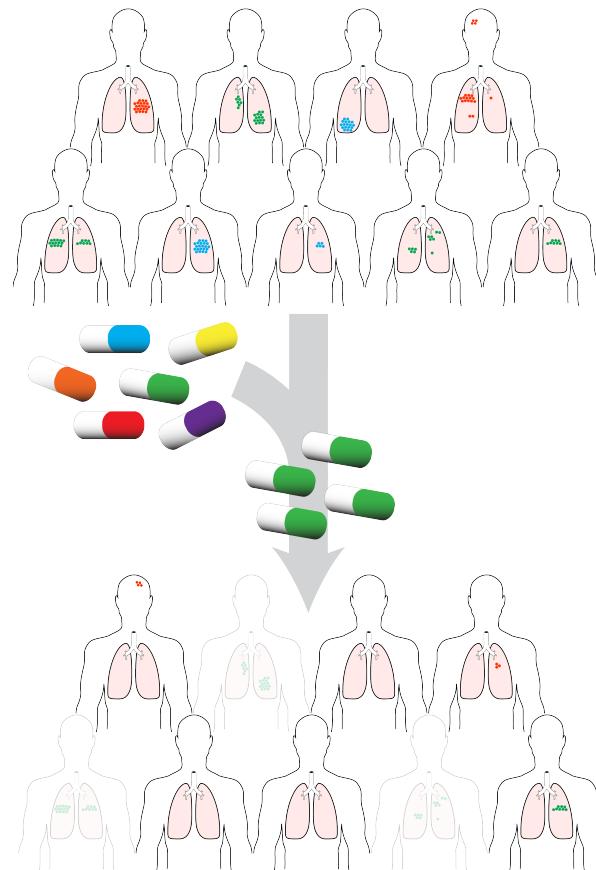
not exhaustive



Precision medicine



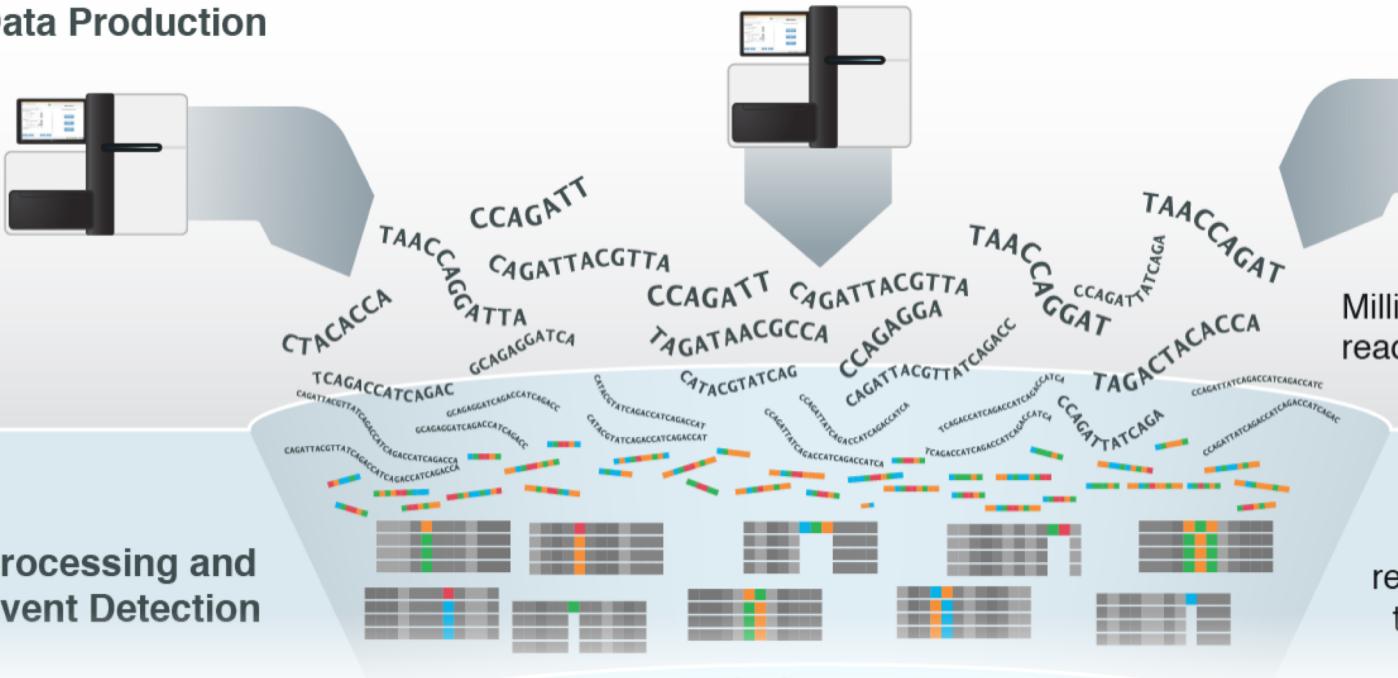
Conventional therapies





Precision Oncology Bottleneck

1. Data Production



Millions of raw sequence reads are produced for a patient tumor.

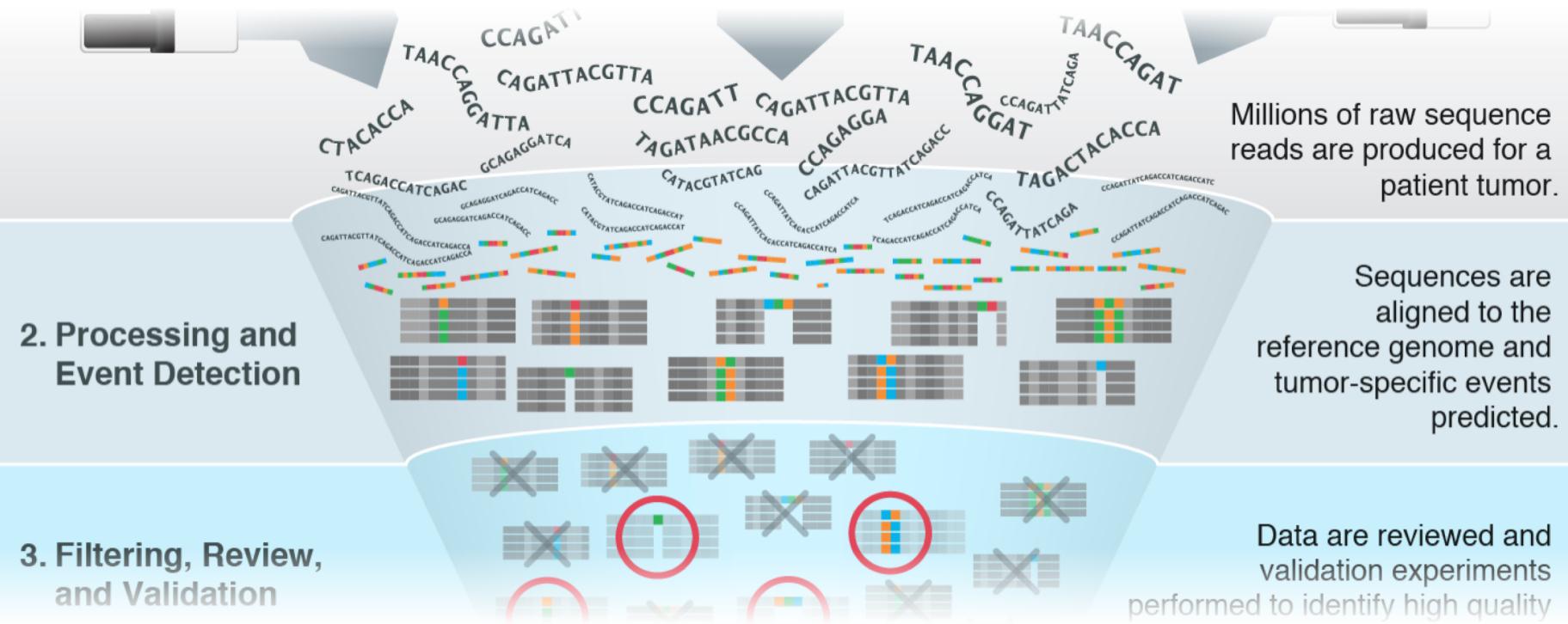
2. Processing and Event Detection

Sequences are aligned to the reference genome and tumor-specific events predicted.





Precision Oncology Bottleneck



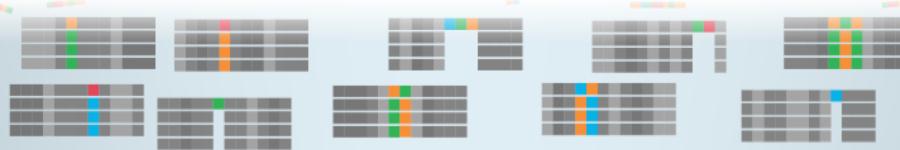
@handlerwagner

Good BM, Ainscough BJ, McMichael JF, Su AI†, Griffith OL†. 2014. Genome Biology. 15(8):438.



Precision Oncology Bottleneck

2. Processing and Event Detection



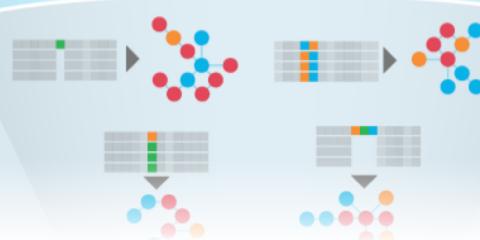
Sequencing data are aligned to the reference genome and tumor-specific events predicted.

3. Filtering, Review, and Validation



Data are reviewed and validation experiments performed to identify high quality events.

4. Annotation and Functional Prediction



Events are annotated and scored in an effort to predict events of functional significance.



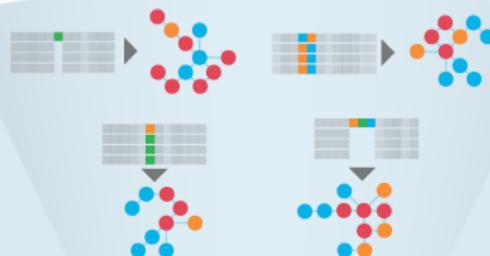
Precision Oncology Bottleneck

3. Filtering, Review, and Validation



Data are reviewed and validation experiments performed to identify high quality events.

4. Annotation and Functional Prediction



Events are annotated and scored in an effort to predict events of functional significance.

5. Interpretation and Report Generation



A genome analyst attempts to interpret, prioritize, and summarize functionally significant events in the context of published literature, clinical trials, and a



Precision Oncology Bottleneck

4. Annotation and Functional Prediction



Events are annotated and scored in an effort to predict events of functional significance.

5. Interpretation and Report Generation



A genome analyst attempts to interpret, prioritize, and summarize functionally significant events in the context of published literature, clinical trials, and a multitude of knowledgebases.

6. Clinical Application

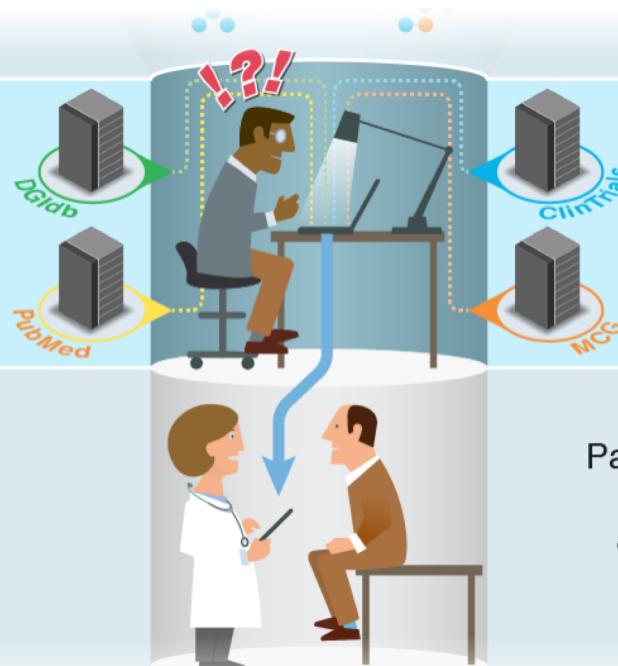


Pathologists and oncologists evaluate the significance of potentially clinically actionable events, and incorporate their research into patient care.



Precision Oncology Bottleneck

5. Interpretation and Report Generation



A genome analyst attempts to interpret, prioritize, and summarize functionally significant events in the context of published literature, clinical trials, and a multitude of knowledgebases.

6. Clinical Application

Pathologists and oncologists evaluate the significance of potentially clinically actionable events, and incorporate their research into patient care.



Precision Oncology Bottleneck





Precision Oncology Bottleneck



How do we alleviate this bottleneck?



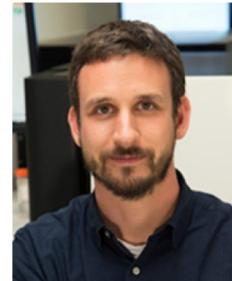
CIViC

CLINICAL INTERPRETATIONS OF
VARIANTS IN CANCER

civicdb.org



Obi Griffith



Malachi Griffith

CIViC principles emphasize collaborative, open sharing of interpretations

Public contributions, open discussion, curation standards and expert review

Researchers, clinicians,
patient advocates and
others

Public domain (CC0)
license

No fees, anonymous access

Content provenance
and creator
acknowledgement

Structured data and
APIs



What is a variant of *significance* in cancer?

A Prognostic of survival change

- *Biallelic CEBPA mutations are associated with improved overall survival in patients with acute myeloid leukemia*

E Predictive of therapeutic response

- *BRAF V600E predicts sensitivity to vemurafenib*

Q Diagnostic of tumor subtype

- *DNAJB1-PRKACA fusion differentiates fibrolamellar hepatocellular carcinoma from conventional HCC*

⚠ Predisposing for cancer development

- *Patients with the RUNX1 Y260* mutation are associated with increased risk of developing acute myeloid leukemia*



What is a variant of *significance* in cancer?



Prognostic of survival change

- Biallelic *CEP1PA* mutations are associated with improved overall survival in patients



Predictive

- *BRAF* V600E



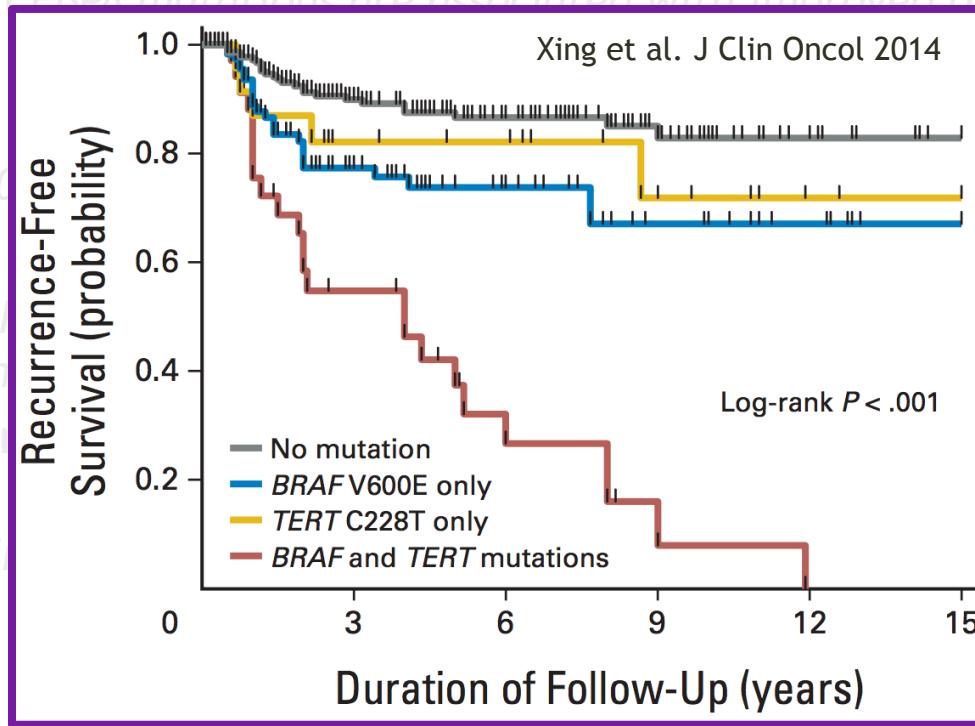
Diagnostic

- *DNAJB1*-like genes from cancer cells

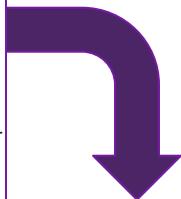
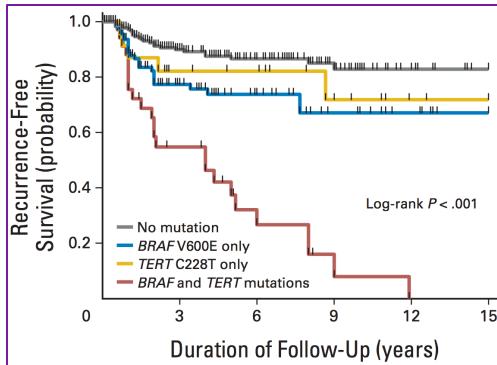


Predisposition

- Patients with increased risk of developing cancer



CIViC provides a literature curation interface for the interpretation of clinically-relevant variants



Structured Evidence Item

EVIDENCE EID656

Evidence Summary Evidence Talk

Submitted by gatoravi Last Modified by kkrysiak Last Reviewed by ahwagner Accepted by NickSpies

In patients with papillary thyroid cancer harboring both BRAF V600E and the TERT promotor mutation C228T (N=35), recurrence-free survival is worse than in patients harboring one of these mutations (N=159 BRAF, N=26 TERT promoter mutated) or no mutations in either gene (N=287)($P<0.001$).

Evidence Level: B - Clinical **Disease:** Papillary Thyroid Carcinoma

Evidence Type: Prognostic **Associated Phenotype:** --

Evidence Direction: Supports **Source:** Xing et al., 2014, J. Clin. Oncol.

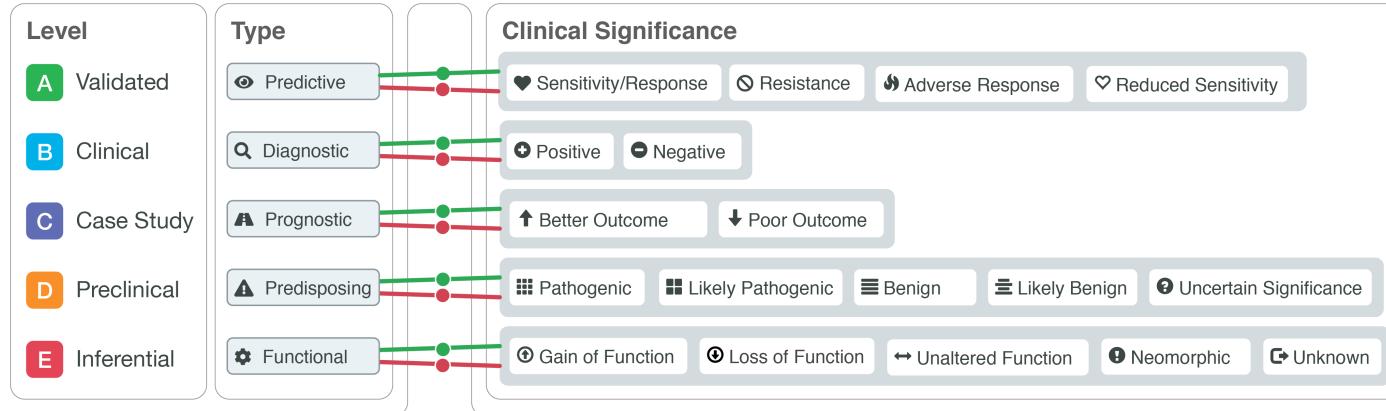
Clinical Significance: Poor Outcome **PubMed ID:** [25024077](#)

Variant Origin: Somatic Mutation **Clinical Trial:** --

Trust Rating: ★★★★★



Curated evidence items are the foundational unit of CIViC



EVIDENCE EID656

[Evidence Summary](#) [Evidence Talk](#) [⚙️](#)

Submitted by gatoravi Last Modified by krysiak Last Reviewed by ahwagner Accepted by NickSpies

In patients with papillary thyroid cancer harboring both BRAF V600E and the TERT promotor mutation C228T (N=35), recurrence-free survival is worse than in patients harboring one of these mutations (N=159 BRAF, N=26 TERT promoter mutated) or no mutations in either gene (N=287)(P<0.001).

Evidence Level: B - Clinical

Evidence Type: Prognostic

Evidence Direction: Supports

Clinical Significance: Poor Outcome

Disease: [Papillary Thyroid Carcinoma](#)

Associated Phenotype: --

Source: Xing et al., 2014, J. Clin. Oncol.

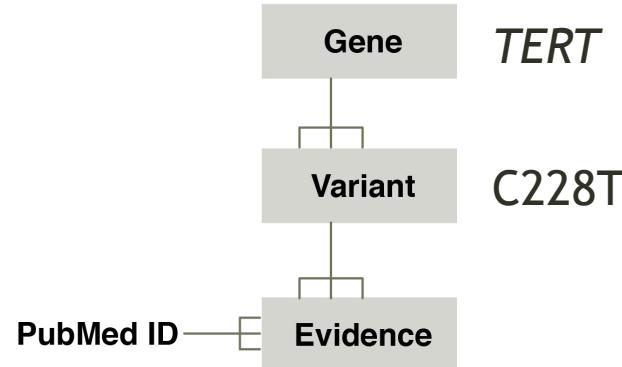
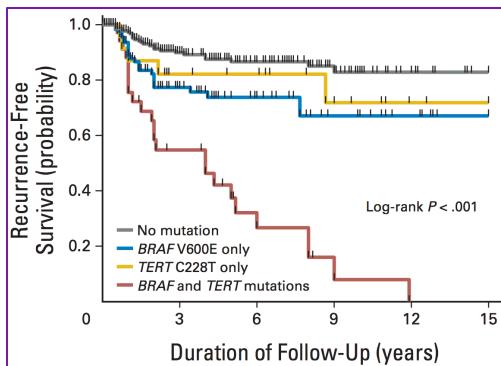
PubMed ID: [25024077](#)

Clinical Trial: --

Trust Rating: ★★★★★



Curated evidence items are the foundational unit of CIViC



EVIDENCE EID656

Submitted by gatoravi Last Modified by kkrysiak Last Reviewed by ahwagner Accepted by NickSpies

In patients with papillary thyroid cancer harboring both BRAF V600E and the TERT promotor mutation C228T (N=35), recurrence-free survival is worse than in patients harboring one of these mutations (N=159 BRAF, N=26 TERT promoter mutated) or no mutations in either gene (N=287)(P<0.001).

Evidence Level: B - Clinical

Evidence Type: Prognostic

Evidence Direction: Supports

Clinical Significance: Poor Outcome

Variant Origin: Somatic Mutation

Disease: Papillary Thyroid Carcinoma

Associated Phenotype: -

Source: Xing et al., 2014, J. Clin. Oncol.

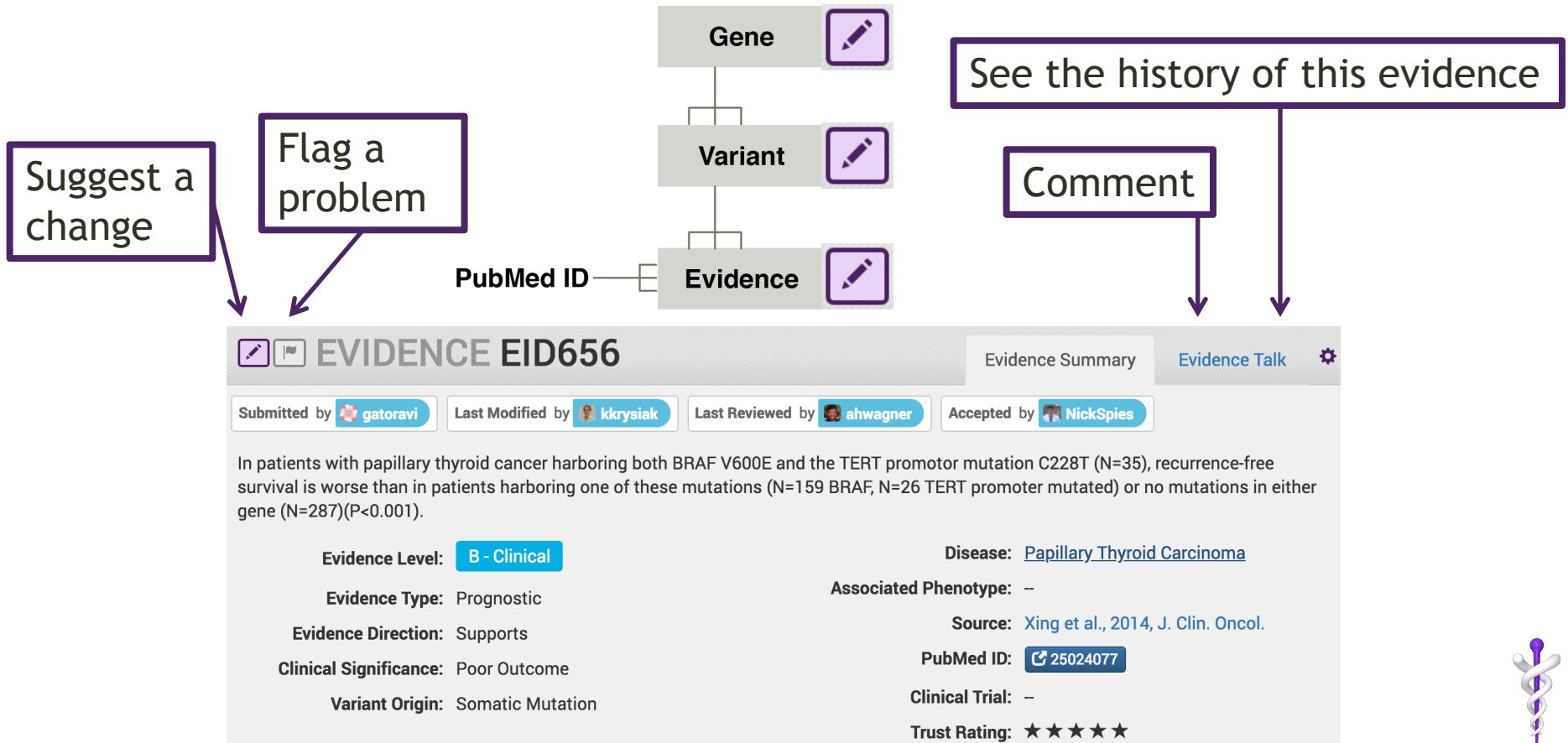
PubMed ID: [25024077](#)

Clinical Trial: -

Trust Rating: ★★★★★



Curated evidence items are the foundational unit of CIViC



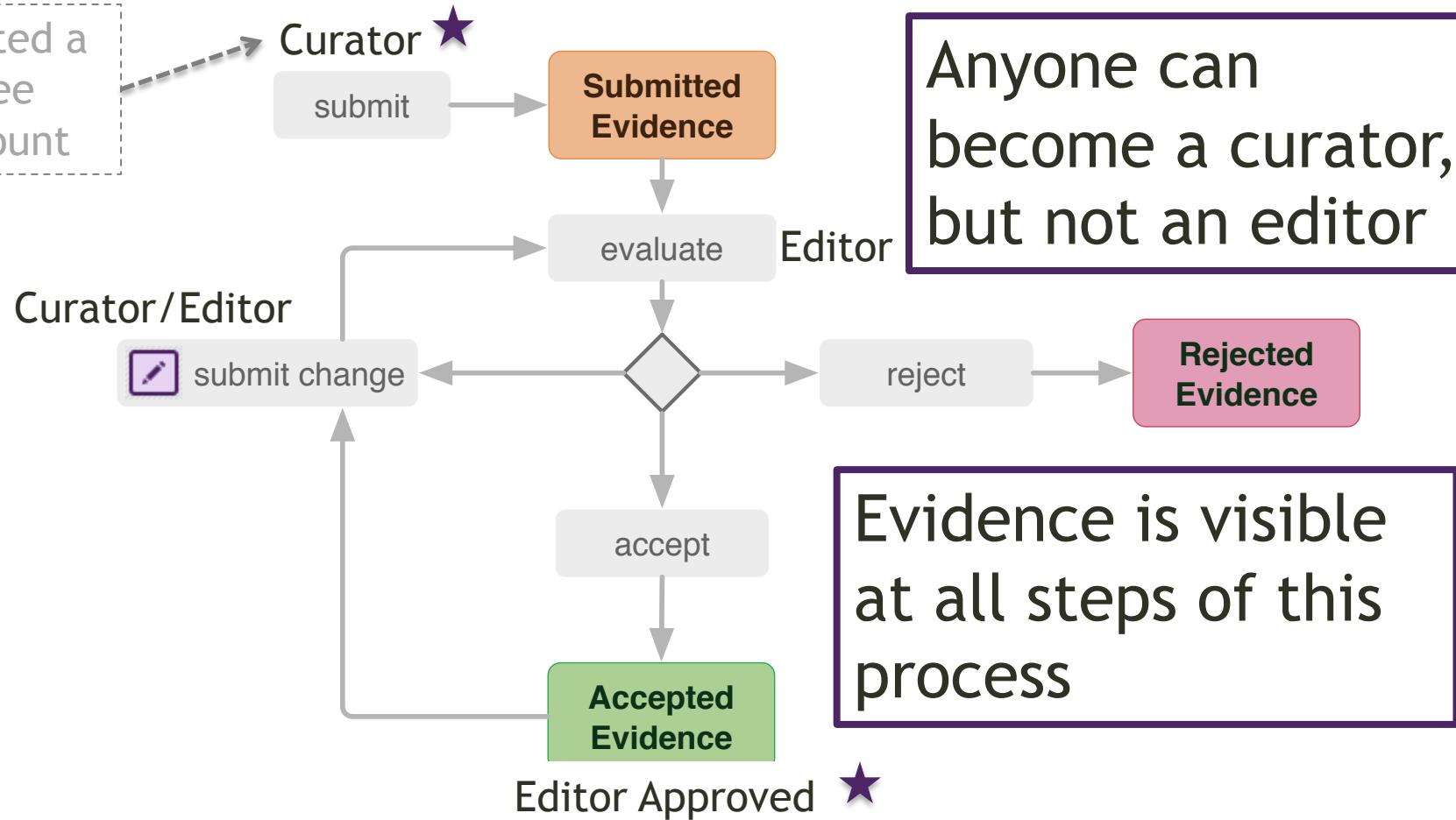
Fundamental requirements of CIViC evidence

- Clinical relevance to cancer
- More than simple observation of a variant in a tumor
- Must not plagiarize the source material
- Published, traceable knowledge (PubMed ID or ASCO abstract)
 - Data available to support assertions
- Must not include personal health information
 - Published peer reviewed case reports are acceptable



The CIViC evidence item life cycle is designed to promote quality, provenance, transparency and adaptation as knowledge evolves

Created a free account



The goal of curating evidence is to document the clinical relevance of genes and variants, and ultimately to support final assertions

CIViC

About Participate Community Help FAQ MyCivc (158)

Go to Genes & Variants Go! BROWSE SEARCH ACTIVITY ADD ▾

GENE ERBB2

Last Modified by obigriffith Last Reviewed by obigriffith

Name: erb-b2 receptor tyrosine kinase 2
Entrez Symbol: ERBB2 **Entrez ID:** 2064

Aliases: CD340, HER-2, HER-2/neu, HER2, MLN 19, NEU, NGL, TKR1

Chromosome: 17 **Start:** 37844167 **End:** 37886679 **Strand:** 1 (GRCh37)

Protein Domains: Furin-like cysteine-rich domain, Furin-like repeat, Growth factor receptor cysteine-rich domain, Growth factor receptor domain 4, Leucine-rich repeat domain, L domain-like... (8 more)

Pathways: miR-targeted genes in muscle cell - TarBase, Leptin signaling pathway, Prolactin Signaling Pathway, Signaling Pathways in Glioblastoma, Integrated Pancreatic Cancer Pathway... (124 more)

[View MyGene.info Details](#)

ERBB2 Variants & Variant Group

Show all: filter variants... Add Group

AMPLIFICATION	D769H	D769Y	DEL 755-759	ERBB2 G776INSV_G/C	G309A	G776L
L753E	L755P	L755S	L755W	L768S	L865M	M774DELINSWLV
M774INSAYVM	MUTATION	N857S	NON-AMPLIFICATION	OVEREXPRESSION	P780INS	
R678Q	R896C	S310F/Y	SERUM LEVELS	T798M	T862A	V773
V773L	V777L	V842I	V772_A775DIP			

VARIANT V777L

Variant Summary Variant Talk

Last Modified by Lymzyell Last Reviewed by alhwagner Last Commented On by alhwagner

Aliases: V777L, V762L, VAL777LEU, and RS121913471 **Allele Registry ID:** CA135387

Sources: Bose et al., 2013, Cancer Discov

HGVIS Expressions: NM_00448.3:c.2329G>T, NP_00449.2:p.Val777Leu, ENST00000269571.5:c.2329G>T, and NC_00017.10:g.37881000G>T

ClinVar ID: 44991 **ClinVar Clinical Significance:** Likely pathogenic

COSMIC ID: COSM14062 **dbSNP RSIID:** rs121913471 **HGVIS ID:** chr17:g.37881000G>T

SnpEff Effect: missense variant **SnpEff Impact:** MODERATE **gnomAD Adj. AF:** –

[View MyVariant.info Details](#)

EVIDENCE EID288

Submitted by NickSpies Accepted by MyCivc

In MCF10A cell lines, the V777L mutation was shown to be sensitive to neratinib.

Evidence Level: D - Preclinical

Evidence Type: Predictive

Evidence Direction: Supports

Clinical Significance: Sensitivity

Variant Origin: Somatic Mutation

Drug: Neratinib

Disease: Breast Cancer

Associated Phenotype: –

Citation: Bose et al., 2013, Cancer Discov

PubMed ID: 23220880

Clinical Trial: –

Trust Rating: ★★★★

Variant Summary Variant Talk

Ref. Build: GRCh37 Ensembl Version: 75

Chr.	Start	Stop	Ref. s	Var. Bases
17	37881000	37881000	G	T

Rep. Transcript: ENST00000269571.5 [Edit Coordinates](#)

ClinVar ID: 44991 **ClinVar Clinical Significance:** Likely pathogenic

COSMIC ID: COSM14062 **dbSNP RSIID:** rs121913471 **HGVIS ID:** chr17:g.37881000G>T

SnpEff Effect: missense variant **SnpEff Impact:** MODERATE **gnomAD Adj. AF:** –

[View MyVariant.info Details](#)

Get Data Help

EID	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR	
5817	In a cohort of 205 Her2-recept...	Her2-receptor Positive Breast ...	Trastuzumab Emtansine, Tras...	C	🕒	🕒	🕒	🕒	🕒	3 ★
288	In MCF10A cell lines, the V777...	Breast Cancer	Neratinib	D	🕒	🕒	🕒	🕒	🕒	5 ★
1177	Colon cancer patient derived ...	Colon Cancer	Lapatinib, Neratinib, Trastuzu...	D	🕒	🕒	🕒	🕒	🕒	4 ★
4452	bp] In an in vitro study, NCI-H...	Colorectal Cancer	Cetuximab	D	🕒	🕒	🕒	🕒	🕒	–

Evidence Summary Evidence Talk

Submitted by NickSpies Accepted by MyCivc

In MCF10A cell lines, the V777L mutation was shown to be sensitive to neratinib.

Evidence Level: D - Preclinical

Evidence Type: Predictive

Evidence Direction: Supports

Clinical Significance: Sensitivity

Variant Origin: Somatic Mutation

Drug: Neratinib

Disease: Breast Cancer

Associated Phenotype: –

Citation: Bose et al., 2013, Cancer Discov

PubMed ID: 23220880

Clinical Trial: –

Trust Rating: ★★★★

Assertions are built from a collection of evidence items

- Each assertion aggregates multiple CIViC evidence items into a single clinical consensus
 - Specific variant, disease, evidence type, [drug]
- Underlying data is carefully organized and transparent, allowing for rapid updates and reassessment as new information emerges
- Once finalized, assertions are submitted to ClinVar



Assertions - example

NSCLC with EGFR
L858R mutation is
sensitive to
erlotinib and
gefitinib

- NCCN guidelines
- FDA approval

...

ASSERTION AID5

Submitted by arpaddanos Last Modified by ebarnell Last Reviewed by arpaddanos Accepted by ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: -

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to erlotinib or gefitinib.

Description: L858R is among the most common sensitizing EGFR mutations in NSCLC, and is assessed via DNA mutational analysis including Sanger sequencing and next generation sequencing methods. Tyrosine kinase inhibitors erlotinib and gefitinib are associated with improved progression free survival over chemotherapy in EGFR L858R patients. NCCN guidelines recommend (category 1) erlotinib and gefitinib for NSCLC with sensitizing EGFR mutations, along with afatinib and osimertinib.

ClinVar ID 16609	ClinVar Clinical Significance Pathogenic	
COSMIC ID COSM6224	dbSNP RSID rs121434568	HGVS ID chr7:g.55259515T>G
SnpEff Effect structural interaction variant	SnpEff Impact HIGH	gnomAD Adj. AF --
View MyVariant.info Details		

Assertion Type: Predictive
Assertion Direction: Supports
Clinical Significance: Sensitivity
Drugs: Erlotinib and Gefitinib
Drug Interaction Type: Substitutes
AMP Category: Tier I - Level A
NCCN Guideline: Non-Small Cell Lung Cancer (v3.2018)
Regulatory Approval:
FDA Companion Test:

Evidence Grid Evidence Cards

Evidence Supporting AID5 14 total items

EID	GENE	VARIANT	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR
2994	EGFR	L858R	On May 14, 2013, the...	Non-small Cell Lung ...	Erlotinib	A					5★
2621	EGFR	L858R	In a phase 3 clinical t...	Non-small Cell Lung ...	Gefitinib	B					4★
1665	EGFR	L858R	90 NSCLC patients w...	Non-small Cell Lung ...	Gefitinib	B					3★
2634	EGFR	L858R	In a phase 3 clinical t...	Lung Adenocarcinoma	Gefitinib	B					3★
229	EGFR	L858R	There is no statistica...	Non-small Cell Lung ...	Erlotinib, Gefitinib (S...	B					3★
885	EGFR	L858R	A randomized phase ...	Non-small Cell Lung ...	Erlotinib	B					3★
2624	EGFR	L858R	Mutational profiling ...	Lung Adenocarcinoma	Gefitinib	C					1★

RULES

Gene and Variant must exist in CIViC

ASSERTION AID5

Submitted by arpaddanos Last Modified by ebarnell Last Reviewed by arpaddanos Accepted by ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: –

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to erlotinib or gefitinib.

Description: L858R is among the most common sensitizing EGFR mutations in NSCLC, and is assessed via DNA mutational analysis including Sanger sequencing and next generation sequencing methods. Tyrosine kinase inhibitors erlotinib and gefitinib are associated with improved progression free survival over chemotherapy in EGFR L858R patients. NCCN guidelines recommend (category 1) erlotinib and gefitinib for NSCLC with sensitizing EGFR mutations, along with afatinib and osimertinib.

ClinVar ID	ClinVar Clinical Significance	
16609	Pathogenic	
COSMIC ID COSM6224	dbSNP RSID rs121434568	HGVS ID chr7:g.55259515T>G
SnpEff Effect structural interaction variant	SnpEff Impact HIGH	gnomAD Adj. AF --
View MyVariant.info Details		

MyVariant.info

Assertion Type: Predictive
Assertion Direction: Supports
Clinical Significance: Sensitivity
Drugs: Erlotinib and Gefitinib
Drug Interaction Type: Substitutes
AMP Category: Tier I - Level A
NCCN Guideline: Non-Small Cell Lung Cancer (v3.2018)
Regulatory Approval: ✓
FDA Companion Test: ✓

Evidence Grid Evidence Cards

Evidence Supporting AID5 14 total items

EID	GENE	VARIANT	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR
2994	EGFR	L858R	On May 14, 2013, the...	Non-small Cell Lung ...	Erlotinib	A					5 ★
2621	EGFR	L858R	In a phase 3 clinical t...	Non-small Cell Lung ...	Gefitinib	B					4 ★
1665	EGFR	L858R	90 NSCLC patients w...	Non-small Cell Lung ...	Gefitinib	B					3 ★
2634	EGFR	L858R	In a phase 3 clinical t...	Lung Adenocarcinoma	Gefitinib	B					3 ★
229	EGFR	L858R	There is no statistica...	Non-small Cell Lung ...	Erlotinib, Gefitinib (S...	B					3 ★
885	EGFR	L858R	A randomized phase ...	Non-small Cell Lung ...	Erlotinib	B					3 ★
2624	EGFR	L858R	Mutational profiling ...	Lung Adenocarcinoma	Gefitinib	C					1 ★

RULES

Gene and Variant must exist in CIViC

Assertion AID5

Submitted by arpaddanos Last Modified by ebarnell Last Reviewed by arpaddanos Accepted by ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: –

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to erlotinib or gefitinib.

Description: L858R is among the most common sensitizing EGFR mutations in NSCLC, and is assessed via DNA mutational analysis including Sanger sequencing and next generation sequencing methods. Tyrosine kinase inhibitors erlotinib and gefitinib are associated with improved progression free survival over chemotherapy in EGFR L858R patients. NCCN guidelines recommend (category 1) erlotinib and gefitinib for NSCLC with sensitizing EGFR mutations, along with afatinib and osimertinib.

ClinVar ID	ClinVar Clinical Significance	
16609	Pathogenic	
COSMIC ID	dbSNP RSID	HGVS ID
COSM6224	rs121434568	chr7:g.55259515T>G
SnpEff Effect	SnpEff Impact	gnomAD Adj. AF
structural interaction variant	HIGH	–

[View MyVariant.info Details](#)

Assertion Type: Predictive
Assertion Direction: Supports
Clinical Significance: Sensitivity
Drugs: Erlotinib and Gefitinib
Drug Interaction Type: Substitutes
AMP Category: Tier I - Level A
NCCN Guideline: Non-Small Cell Lung Cancer (v3.2018)
Regulatory Approval: ✓
FDA Companion Test: ✓

Evidence Grid Evidence Cards

Evidence Supporting AID5 14 total items

EID	GENE	VARIANT	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR
2994	EGFR	L858R	On May 14, 2013, the...	Non-small Cell Lung ...	Erlotinib	A	🕒	👉	❤️	...	5 ★
2621	EGFR	L858R	In a phase 3 clinical t...	Non-small Cell Lung ...	Gefitinib	B	🕒	👉	❤️	...	4 ★
1665	EGFR	L858R	90 NSCLC patients w...	Non-small Cell Lung ...	Gefitinib	B	🕒	👉	❤️	...	3 ★
2634	EGFR	L858R	In a phase 3 clinical t...	Lung Adenocarcinoma	Gefitinib	B	🕒	👉	❤️	...	3 ★
229	EGFR	L858R	There is no statistica...	Non-small Cell Lung ...	Erlotinib, Gefitinib (S...	B	🕒	👉	❤️	...	3 ★
885	EGFR	L858R	A randomized phase ...	Non-small Cell Lung ...	Erlotinib	B	🕒	👉	❤️	...	3 ★
2624	EGFR	L858R	Mutational profiling ...	Lung Adenocarcinoma	Gefitinib	C	🕒	👉	❤️	...	1 ★

Evidence must be fully curated

RULES

ASSERTION AID5

[Assertion Summary](#)[Assertion Talk](#)Submitted by  arpaddanosLast Modified by  ebarnellLast Reviewed by  arpaddanosAccepted by  ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: -

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to erlotinib or gefitinib.

Description: L858R is among the most common sensitizing EGFR mutations in NSCLC, and is assessed via DNA mutational analysis including Sanger sequencing and next generation sequencing methods. Tyrosine kinase inhibitors erlotinib and gefitinib are associated with improved progression free survival over chemotherapy in EGFR L858R patients. NCCN guidelines recommend (category 1) erlotinib and gefitinib for NSCLC with sensitizing EGFR mutations, along with afatinib and osimertinib.

ClinVar ID	ClinVar Clinical Significance	
16609	Pathogenic	
COSMIC ID	dbSNP RSID	HGVs ID
COSM6224	rs121434568	chr7:g.55259515T>G
SnpEff Effect	SnpEff Impact	gnomAD Adj. AF
structural interaction variant	HIGH	—
View MyVariant.info Details		

Assertion Type: Predictive

Assertion Direction: Supports

Clinical Significance: Sensitivity

Drugs: Erlotinib and Gefitinib

Drug Interaction Type: Substitutes

AMP Category: Tier I - Level A

NCCN Guideline: Non-Small Cell Lung Cancer (v3.2018)

Regulatory Approval: ✓

FDA Companion Test: ✓

Gene and Variant must exist in CIViC

Critical rationale / decisions are described

Evidence must be fully curated

[Evidence Grid](#) [Evidence Cards](#)

Evidence Supporting AID5 14 total items

EID	GENE	VARIANT	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR	
2994	EGFR	L858R	On May 14, 2013, the...	Non-small Cell Lung ...	Erlotinib	A	🕒	👉	❤️	...	5 ★	
2621	EGFR	L858R	In a phase 3 clinical t...	Non-small Cell Lung ...	Gefitinib	B	🕒	👉	❤️	...	4 ★	
1665	EGFR	L858R	90 NSCLC patients w...	Non-small Cell Lung ...	Gefitinib	B	🕒	👉	❤️	...	3 ★	
2634	EGFR	L858R	In a phase 3 clinical t...	Lung Adenocarcinoma	Gefitinib	B	🕒	👉	❤️	...	3 ★	
229	EGFR	L858R	There is no statistica...	Non-small Cell Lung ...	Erlotinib, Gefitinib (S...	B	🕒	👉	❤️	...	3 ★	
885	EGFR	L858R	A randomized phase ...	Non-small Cell Lung ...	Erlotinib	B	🕒	👉	❤️	...	3 ★	
2624	EGFR	L858R	Mutational profiling ...	Lung Adenocarcinoma	Gefitinib	C	🕒	👉	❤️	...	1 ★	

RULES

Gene and Variant must exist in CIViC

Critical rationale / decisions are described

Evidence must be fully curated

ASSERTION AID5

Submitted by arpaddanos Last Modified by ebarnell Last Reviewed by arpaddanos Accepted by ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: –

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to erlotinib or gefitinib.

Description: L858R is among the most common sensitizing EGFR mutations in NSCLC, and is assessed via DNA mutational analysis including Sanger sequencing and next generation sequencing methods. Tyrosine kinase inhibitors erlotinib and gefitinib are associated with improved progression free survival over chemotherapy in EGFR L858R patients. NCCN guidelines recommend (category 1) erlotinib and gefitinib for NSCLC with sensitizing EGFR mutations, along with afatinib and osimertinib.

ClinVar ID 16609	ClinVar Clinical Significance Pathogenic		MyVariant.info
COSMIC ID COSM6224	dbSNP RSID rs121434568	HGVS ID chr7:g.55259515T>G	
SnpEff Effect structural interaction variant	SnpEff Impact HIGH	gnomAD Adj. AF --	

[View MyVariant.info Details](#)

Evidence Grid **Evidence Cards**

Evidence Supporting AID5 14 total items

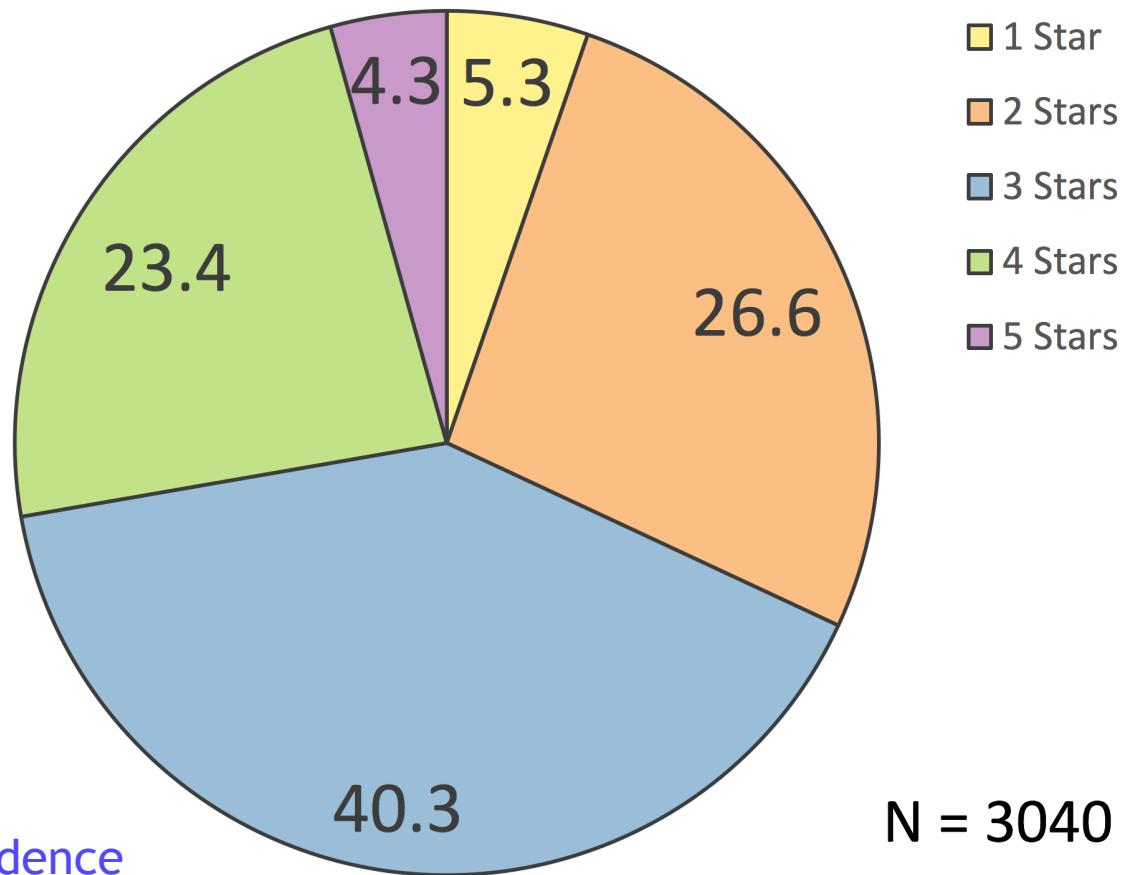
EID	GENE	VARIANT	DESC	DIS	DRUGS	EL	ET	ED	CS	VO	TR
2994	EGFR	L858R	On May 14, 2013, the...	Non-small Cell Lung ...	Erlotinib	A					5★
2621	EGFR	L858R	In a phase 3 clinical t...	Non-small Cell Lung ...	Gefitinib	B					4★
1665	EGFR	L858R	90 NSCLC patients w...	Non-small Cell Lung ...	Gefitinib	B					3★
2634	EGFR	L858R	In a phase 3 clinical t...	Lung Adenocarcinoma	Gefitinib	B					3★
229	EGFR	L858R	There is no statistica...	Non-small Cell Lung ...	Erlotinib, Gefitinib (S...	B					3★
885	EGFR	L858R	A randomized phase ...	Non-small Cell Lung ...	Erlotinib	B					3★
2624	EGFR	L858R	Mutational profiling ...	Lung Adenocarcinoma	Gefitinib	C					1★

Assertion Type: Predictive
Assertion Direction: Supports
Clinical Significance: Sensitivity
Drugs: Erlotinib and Gefitinib
Drug Interaction Type: Substitutes
AMP Category: Tier I - Level A
NCCN Guideline: Non-Small Cell Lung Cancer (v3.2018)
Regulatory Approval:
FDA Companion Test:

At least 1 Evidence item with 3+ stars

Star ratings help guide Assertion creation

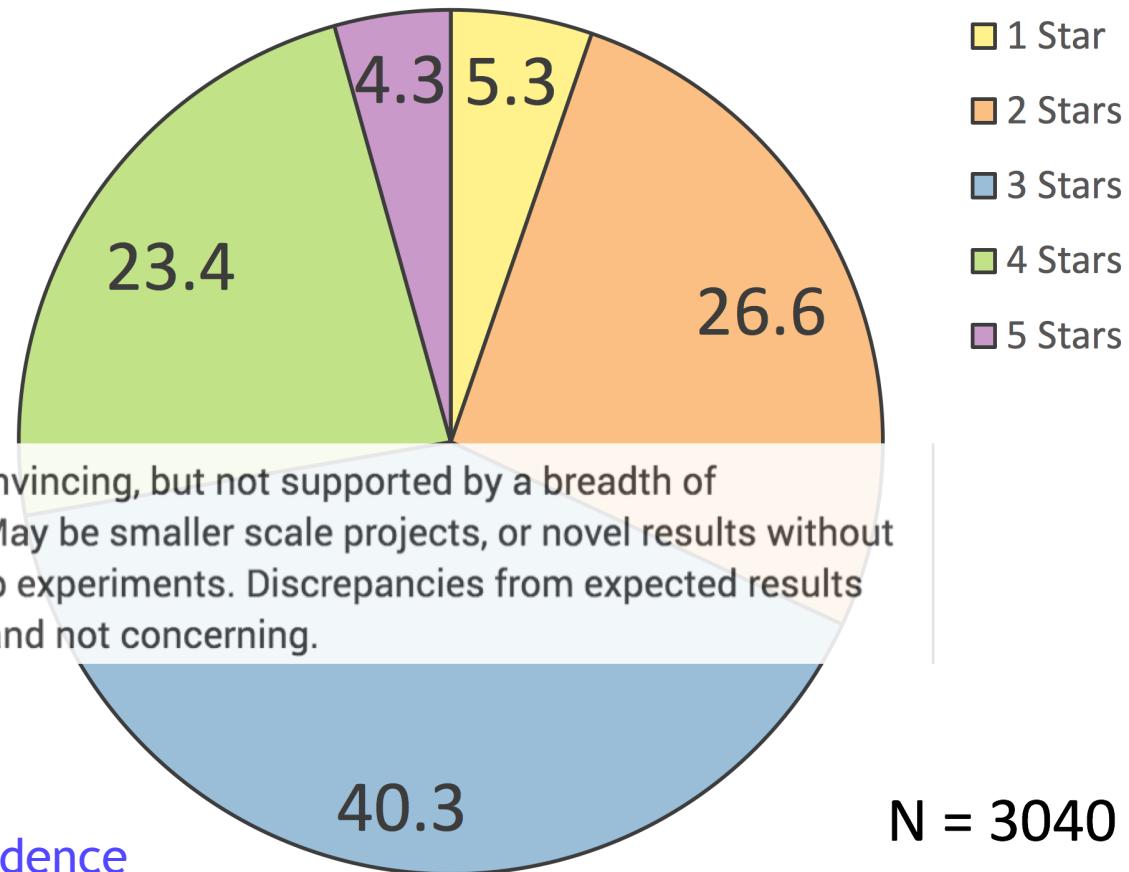
- Moderated evidence is appropriately distributed



<https://civicdb.org/statistics/evidence>

Star ratings help guide Assertion creation

- Moderated evidence is appropriately distributed



<https://civicdb.org/statistics/evidence>

RULES

Assertion AID5

Assertion Summary

Assertion Talk

Submitted by arpaddanos

Last Modified by ebarnell

Last Reviewed by arpaddanos

Accepted by ebarnell

Gene: EGFR Variant: L858R Variant Allele Registry ID: CA126713

Disease: Non-small Cell Lung Carcinoma

Associated Phenotype: -

Summary: Non-small cell lung cancer with EGFR L858R mutation is sensitive to

Assertion Type: Predictive

Assertion Direction: Supports

Tier I: Variants of Strong Clinical Significance

Therapeutic, prognostic & diagnostic

Level A Evidence

FDA-approved therapy
Included in professional guidelines

Level B Evidence

Well-powered studies with consensus from experts in the field

Tier II: Variants of Potential Clinical Significance

Therapeutic, prognostic & diagnostic

Level C Evidence

FDA-approved therapies for different tumor types or investigational therapies
Multiple small published studies with some consensus

Level D Evidence

Preclinical trials or a few case reports without consensus

Tier III: Variants of Unknown Clinical Significance

Not observed at a significant allele frequency in the general or specific subpopulation databases, or pan-cancer or tumor-specific variant databases

No convincing published evidence of cancer association

Tier IV: Benign or Likely Benign Variants

Observed at significant allele frequency in the general or specific subpopulation databases
No existing published evidence of cancer association

8)

Help

TR

▼

5★
4★
3★
3★
3★
3★
1★

How do we scale up without compromising quality (ClinGen Somatic)?

- Collaboration with the ClinGen Somatic Working Group is critical to achieving broad engagement by domain experts
 - Subha Madhavan and Shruti Rao
- Focused working groups forming (disease, gene, variant type, etc.)
- >40 members
- >250 evidence items submitted

CIViC

About Participate Community Help FAQ MalachiGriffith 13 ▾

Go to Genes & Variants Go! BROWSE SEARCH ACTIVITY ADD ▾

Organization Summary

Pediatric Cancer Task Force at ClinGen [Organization Website](#)

Pediatric Cancer Task Force within the ClinGen Somatic Work Group focuses on curation of clinically relevant somatic variants in pediatric cancers.

 ClinGen Clinical Genome Resource

Members

 <p>Name Deb Irene Role Curator Expertise Research Scientist DebIrene</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>	 <p>Name Shruti Rao Role Curator Expertise -- ShrutiRao</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>	 <p>Name Chimene Kesserwan Role Curator Expertise -- ChimeneKesserwan</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID 0000-0001-6043-2065</p>	 <p>Name Laura Corson Role Curator Expertise Research Scientist LauraCorson</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>
 <p>Name Gordana Raca Role Curator Expertise Clinical Scientist graca</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>	 <p>Name Angshumoy Roy Role Curator Expertise Clinical Scientist AngshumoyRoy</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID 0000-0001-7248-8576</p>	 <p>Name Kristen Lipscomb Sund Role Curator Expertise -- KristenLipscombSund</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>	 <p>Name Nishant Tiwari Role Curator Expertise Clinical Scientist Nishant.Tiwari</p> <p>Organization Pediatric Cancer Task Force at ClinGen ORCID ID --</p>



Creation of standards, SOPs, and training materials is needed

BMC Part of Springer Nature

Genome Medicine

Home About Articles Submission Guidelines

Correspondence | Open Access | Published: 29 November 2019

Standard operating procedure for curation and clinical interpretation of variants in cancer

Arpad M. Danos, Kilannin Krysiak, Erica K. Barnell, Adam C. Coffman, Joshua F. McMichael, Susanna Kiwala, Nicholas C. Spies, Lana M. Sheta, Shahil P. Pema, Lynzey Kujan, Kaitlin A. Clark, Amber Z. Wollam, Shruti Rao, Deborah I. Ritter, Dmitriy Sonkin, Gordana Raca, Wan-Hsin Lin, Cameron J. Grisdale, Raymond H. Kim, Alex H. Wagner, Subha Madhavan, Malachi Griffith & Obi L. Griffith

Genome Medicine 11, Article number: 76 (2019) | Cite this article

380 Accesses | 4 Altmetric | Metrics

Abstract

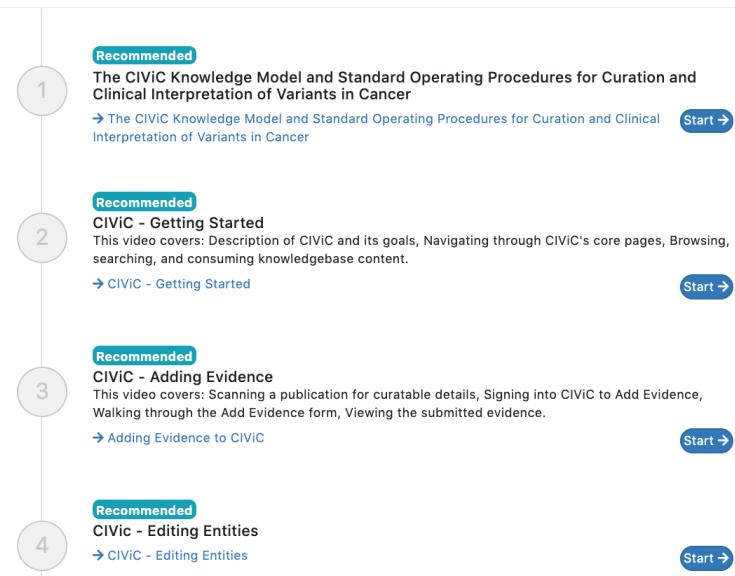
Manually curated variant knowledgebases and their associated knowledge models are serving an increasingly important role in distributing and interpreting variants in cancer. These knowledgebases vary in their level of public accessibility,

Somatic Training Materials

Somatic Variant Training Materials Documents

Interested in Somatic Variant Curation? In order to get involved with our activities, please fill out our volunteer survey: <http://bit.ly/clingenvolunteersurvey>. For questions about existing materials or requests for new materials, contact us at clingen@clinicalgenome.org.

Training Modules Additional Supporting Materials

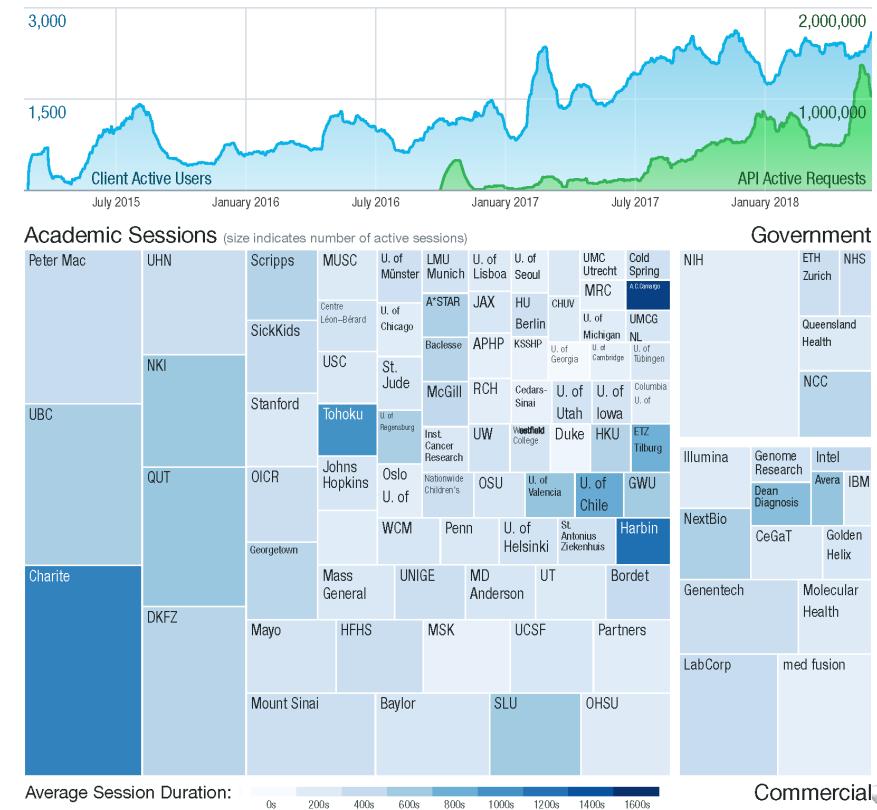
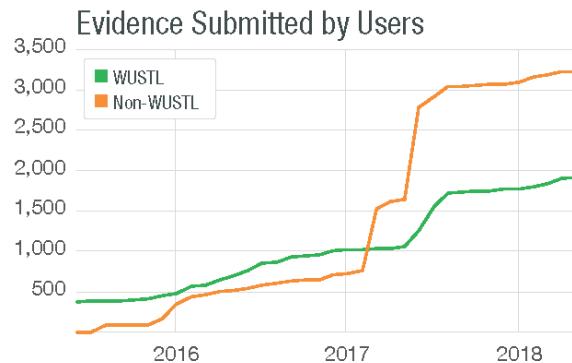


<https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-019-0687-x>

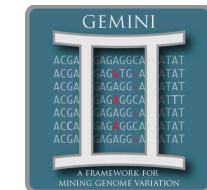
<https://clinicalgenome.org/curation-activities/somatic/training-materials/>

How do we measure success?

- 6,468 evidence lines curated for 2,357 variants, 402 genes, 274 cancer types, 2,374 papers. 24 assertions to date.
- >3,000 users per month
- 191 contributors to date
- >750,000 API requests per month



Open model promotes adoption by commercial and academic data clients



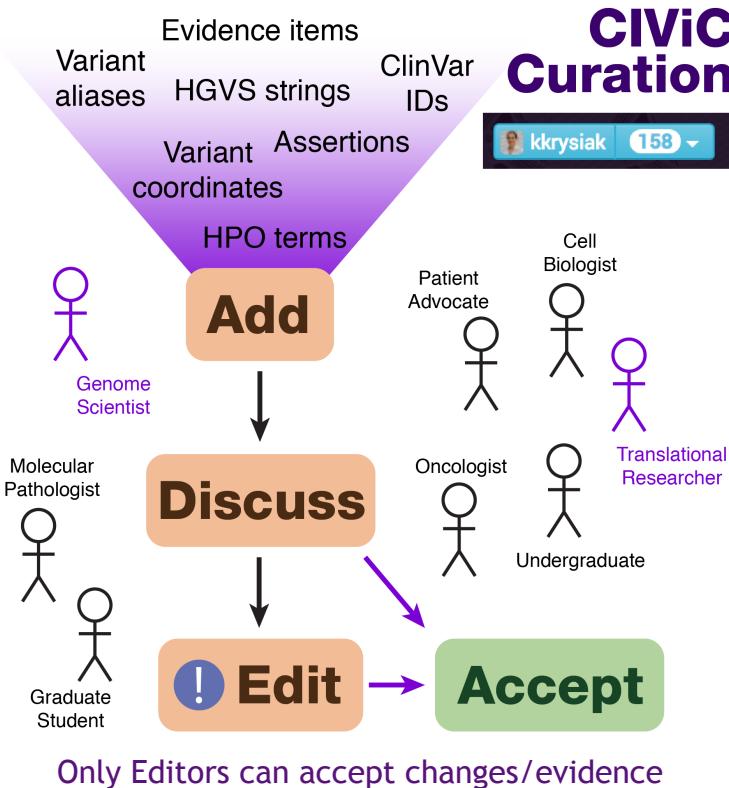
Alissa Clinical Informatics Platform
Alissa Interpret



CIViC has >35 known data clients (<https://civicdb.org/about#data-clients>)



CIViC's open access and crowd-sourcing model has advantages and disadvantages



- **Advantages**
- Scalable
- Promotes transparent consensus building (with provenance)
- Diverse community (remote, asynchronous, global contributions)
- **Disadvantages**
- Quality (perceived and real)
- Difficult to focus curation effort
- Funding





Problem: CIViC is a Silo in a Diverse Knowledge Ecosystem



Variant interpretation knowledge remains siloed

Established Interpretation Knowledgebases (somatic)

- [CIViC \(WashU\)](#)
- [OncoKB \(MSKCC\)](#)
- [JAX-Clinical Knowledgebase \(Jackson lab\)](#)
- [Cancer Genome Interpreter \(Barcelona\)](#)
- [MyCancerGenome \(Vanderbilt\)](#)
- [PMKB \(Cornell\)](#)
- [KnowledgeBase for Precision Oncology \(MD Anderson\)](#)
- [CanDL \(Ohio State\)](#)
- [COSMIC \(Sanger\)](#)

Germline

- [PharmGKB](#)
- [ClinVar](#)
- [ClinGen Evidence Repository](#)

Additionally...

- Many ad hoc “databases” (academic centers / hospitals)
- Many industry examples

Collaborative standards and interoperability work is needed: Variant Interpretation for Cancer Consortium (VICC) - GA4GH Driver Project



...

cancervariants.org



Global Alliance
for Genomics & Health

Variant Interpretation for Cancer

- Gene
- Variant
- Cancer subtype
- Clinical implication: drug sensitivity, drug resistance, adverse response, diagnostic, or prognostic
- Source (e.g., PubMed identifier)
- Curation group

ga4gh.org

Goals/Principles:

- Global integration of clinical cancer variant interpretation
- Standards and guidelines
- Open content for sharing
- Facilitate cross-knowledgebase queries

search.cancervariants.org

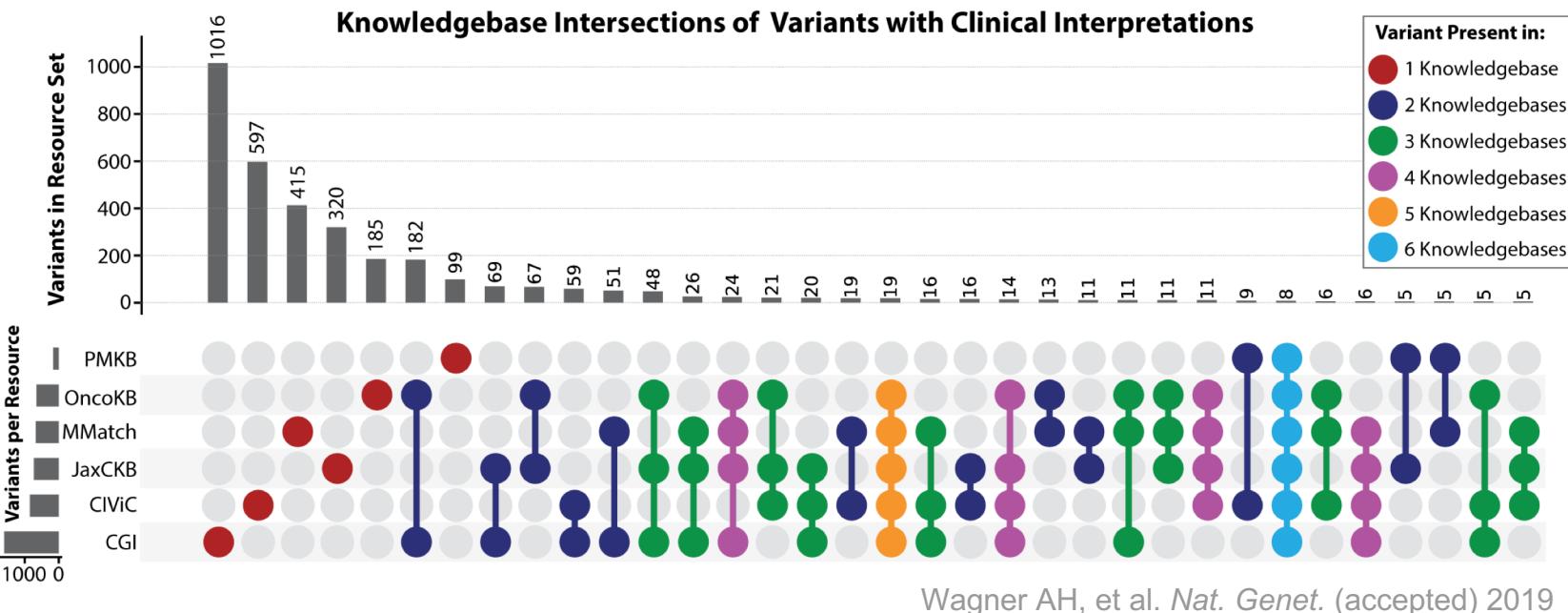




Variant Interpretations are **Heterogeneous**

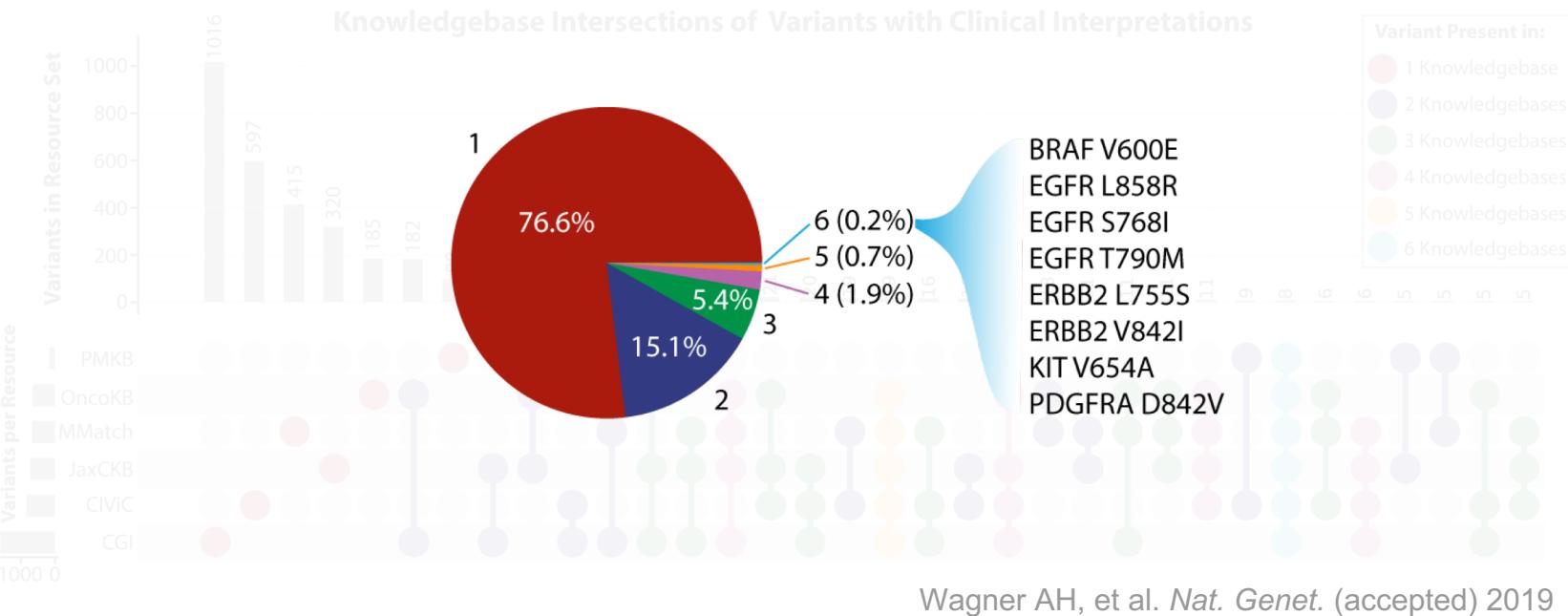


Variant Representation Across Knowledgebases





Variant Representation Across Knowledgebases





Poor overlap of variants creates
inconsistent clinical interpretations

This can be addressed by
aggregating curated knowledge





Variant Interpretations are Structurally Disparate

Diversity in Structure of Variant Interpretations



Gene	BRAF	BRAF (Entrez ID: 673)	BRAF
Isoform	ENST00000288602 / NM_004333.4	ENST00000288602.6	ENST00000288602
Variant	V600E (?????)	V600E (chr7:g.140453136A>T)	V600E (7:140453136:140453136)
Disease	Melanoma	Skin Melanoma (DOID:8923)	Tumor: Melanoma / Tissue: Skin
Drug	Dabrafenib	Dabrafenib + Trametinib	?
Clinical Significance	Known Effect: Sensitive	Supports Sensitivity	?
Evidence Level	2B	A - Validated	Tier 1
Statement	Approved Indications: Dabrafenib is FDA-approved for BRAF V600E mutant unresectable or metastatic melanoma	Open-label, randomized phase 3 trial with 704 patients with metastatic melanoma with a BRAF V600 mutation. Patients were randomized Various B-Raf inhibitors (Vemurafenib, Dabrafenib) have been FDA approved for melanoma therapy in certain settings.



Poor overlap of variants creates
inconsistent clinical interpretations

We can aggregate, but...

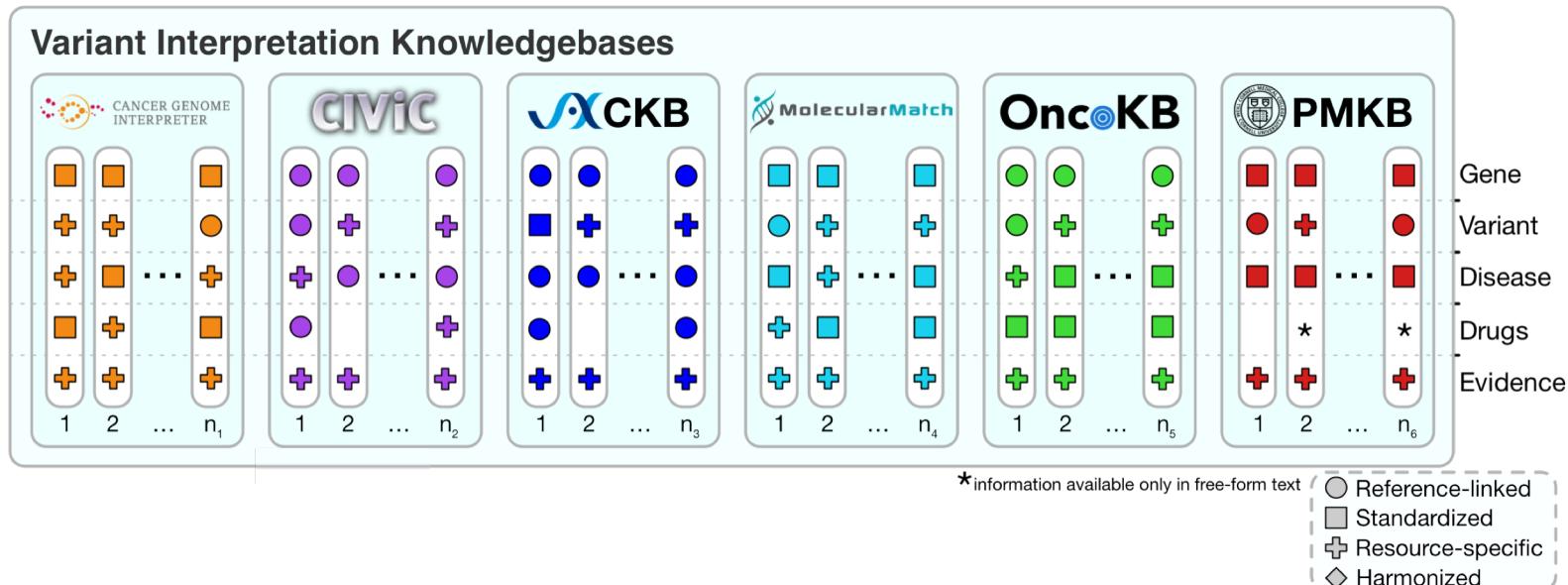
Clinical variant interpretation
knowledgebases are not interoperable

This can be addressed by harmonizing interpretations





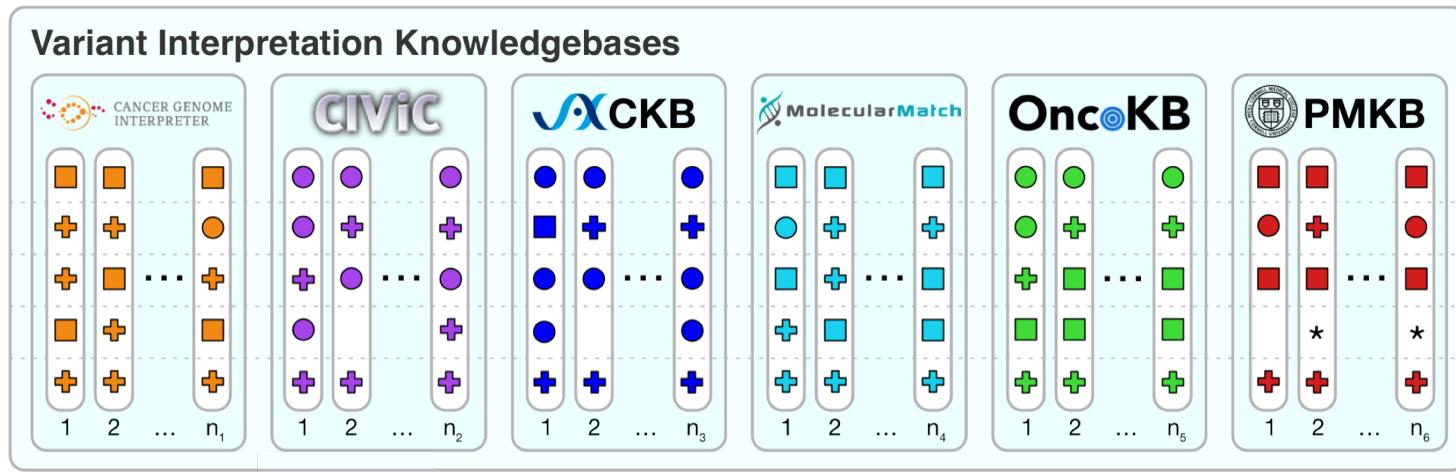
Structure of a Variant Interpretation



Wagner AH, et al. *Nat. Genet.* (accepted) 2019



Structure of a Variant Interpretation

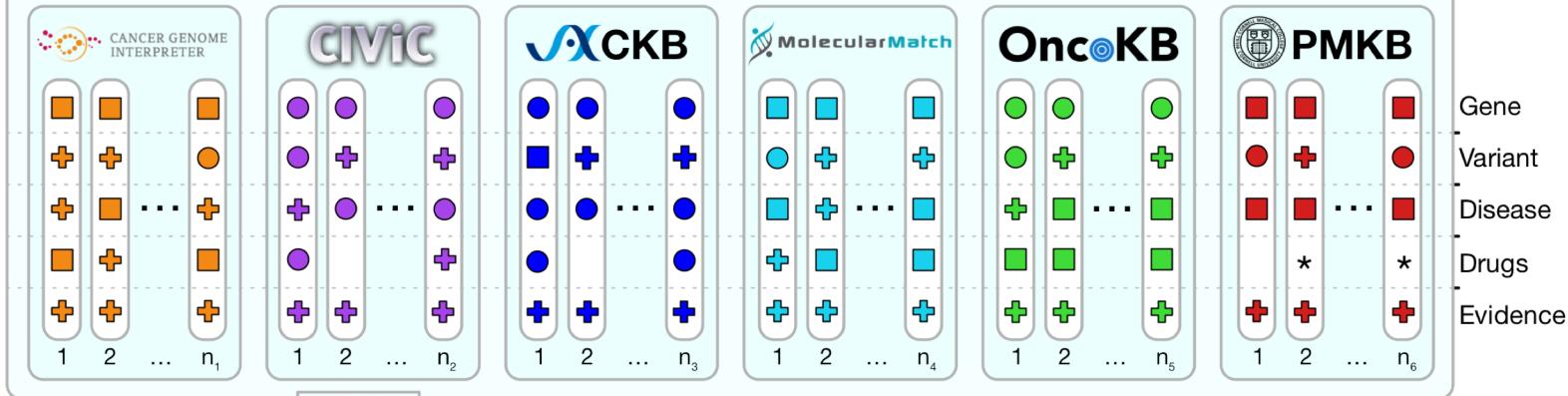


Wagner AH, et al. *Nat. Genet.* (accepted) 2019

AMP/ASCO/CAP Variant Evidence Guidelines							
Evidence Level	Defining Characteristics	CIViC	OncoKB	JAX-CKB	CGI	MMatch	PMKB
Level A (Tier I)	<i>Evidence from professional guidelines or FDA-approved therapies relating to a biomarker and disease.</i>	Level A	Level 1 / 2A /R1	Guideline / FDA Approved	Clinical Practice	Level 1A	Tier 1
Level B (Tier I)	<i>Evidence from clinical trials or other well-powered studies in clinical populations, with expert consensus.</i>	Level B	Level 3A	Phase III	Clinical Trials III-IV	Level 1B	
Level C (Tier II)	<i>Evidence for therapeutic predictive markers from case studies, or other biomarkers from several small studies. Also evidence for biomarker therapeutic predictions for established drugs for different indications.</i>	Predictive Level C	Level 2B, Level 3B	Clinical Study/ Phase I / Phase II	Clinical Trials I-II, Case Reports	Level 2C	Tier 2
Level D (Tier II)	<i>Preclinical findings or case studies of prognostic or diagnostic biomarkers. Also includes indirect findings.</i>	Non-predictive Level C / Level D / Level E	Level 4	Phase 0, Pre-clinical	Pre-clinical Data	Level 2D	

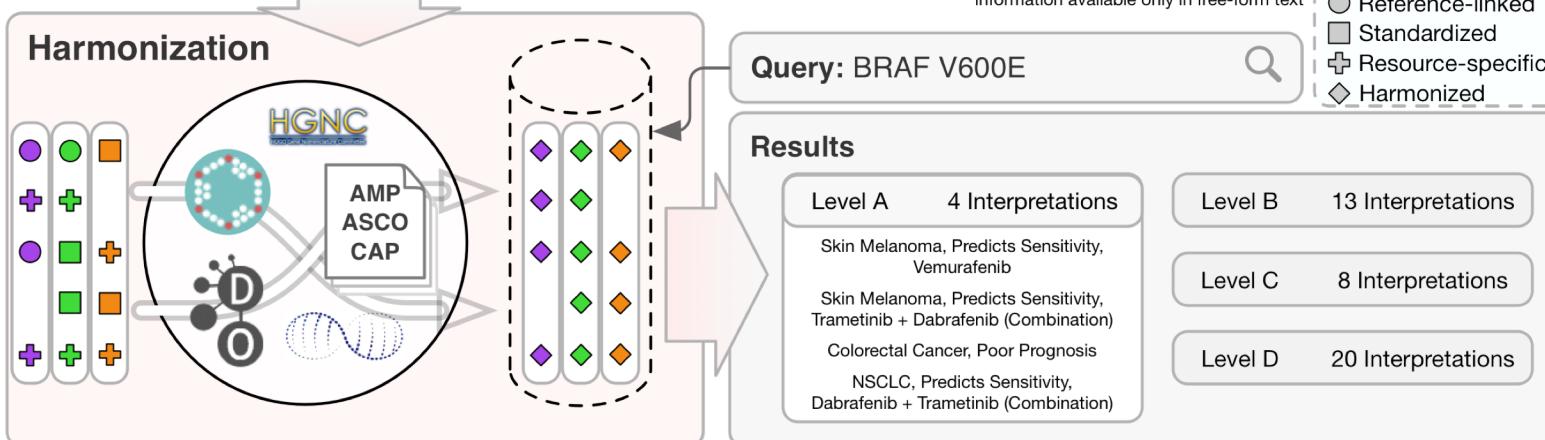


Variant Interpretation Knowledgebases



*information available only in free-form text

- Reference-linked
- Standardized
- + Resource-specific
- ◊ Harmonized



Wagner AH, et al. *Nat. Genet.* (accepted) 2019



Harmonizing Variants: An Expressive and Computable specification for Variation Representation (VR-spec / VR)

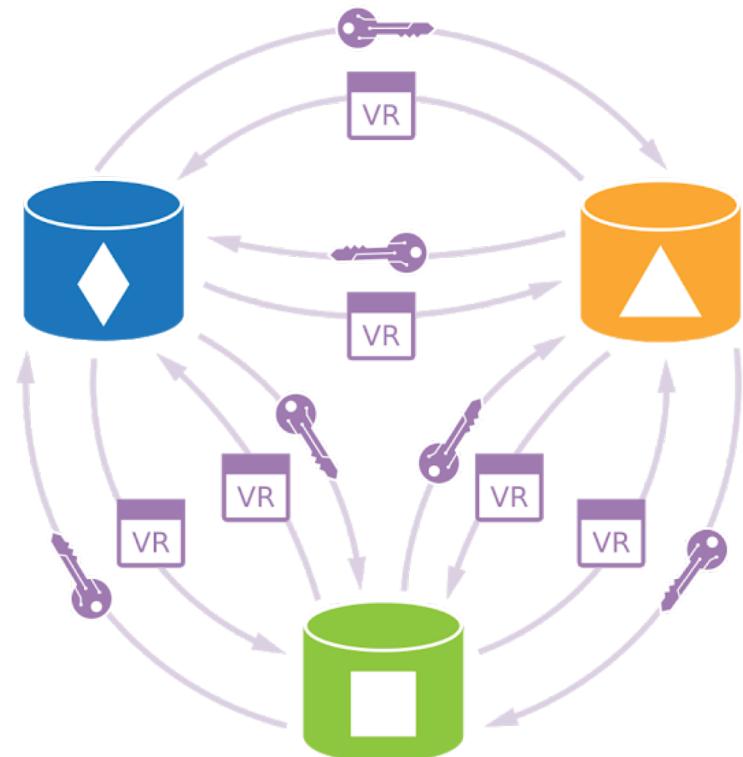


Overview

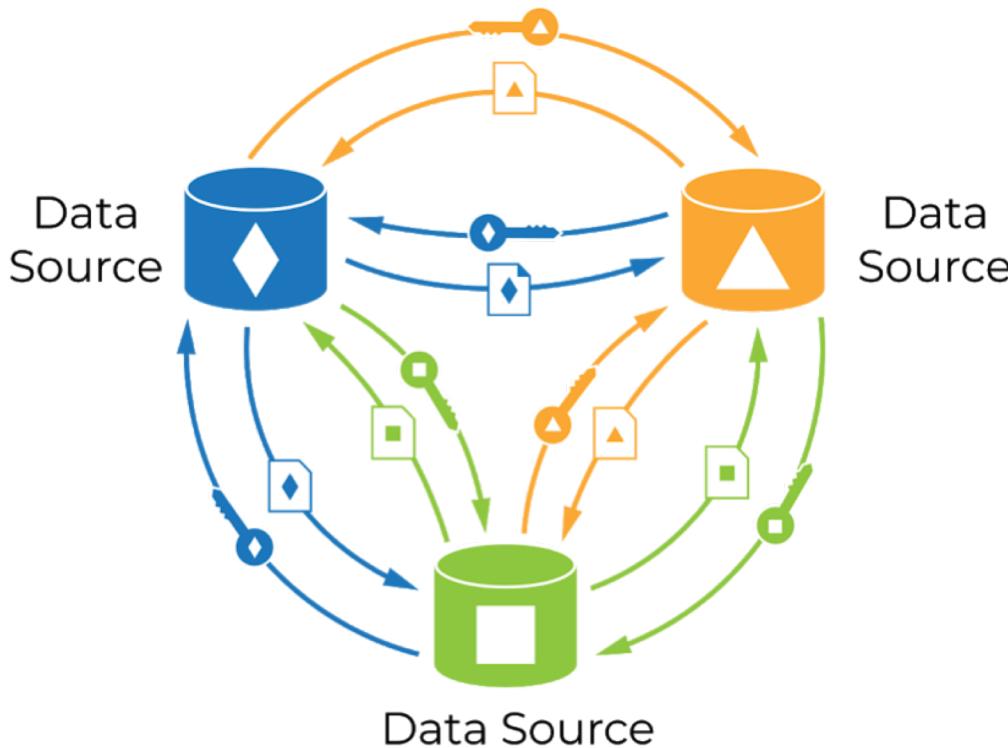
VR aims to standardize the definition, sharing, and identification of biological variation within and between systems.

VR 1.0 Aims:

- An extensible schema for variation
- Minimization of representational ambiguity
- Shared identifiers for variation
- Example Implementation

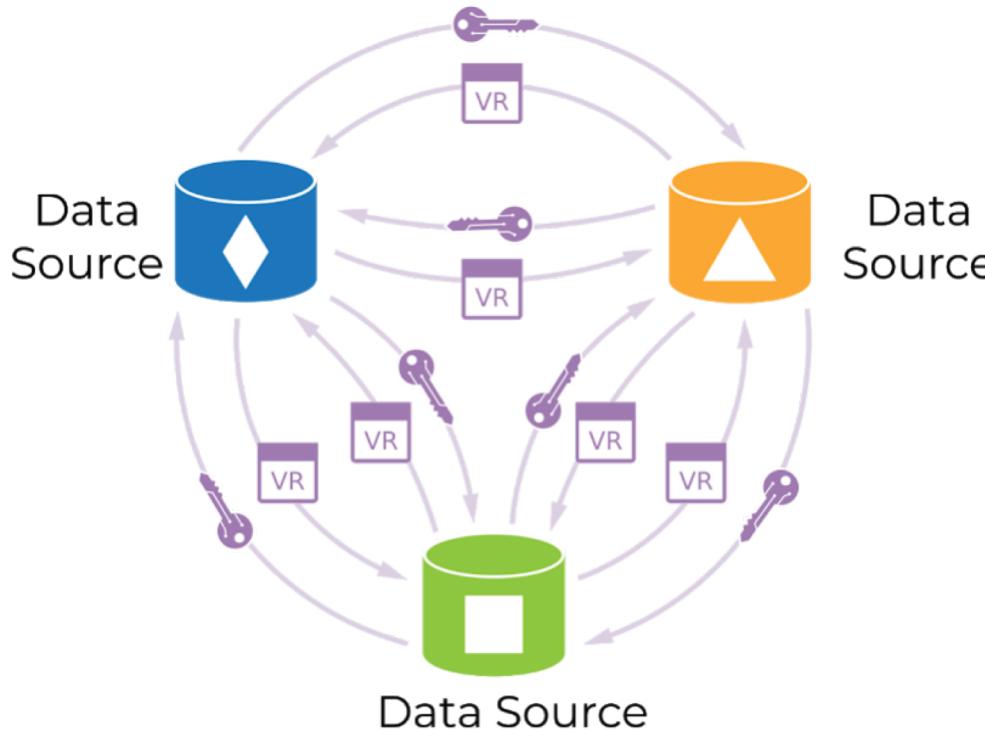


CURRENTLY...



PAIRS OF SYSTEMS COORDINATE KEYS AND FORMATS IN ORDER TO SHARE VARIATION DATA.
ADDING A NEW SYSTEM IS DIFFICULT.

WITH THE VR SPECIFICATION...



SYSTEMS USE A COMMON IDENTIFIER, COMPUTED FROM THE DATA ITSELF, AND A COMMON DATA FORMAT. ADDING A NEW SYSTEM IS MUCH EASIER.

Presentation ≠ Representation

NM_080588.2:c.139_140insC
ENST00000367279:c.139_140insC
NM_080588.2:c.139dup
ENST00000367279:c.139dup



```
{  
  "sequence_id": "ga4gh:SQ.0a1b2c3d",  
  "location": (139,139),  
  "state": "C"  
}
```

Multiple human names
due to choice of accession, normalization,
ins/dup, etc.



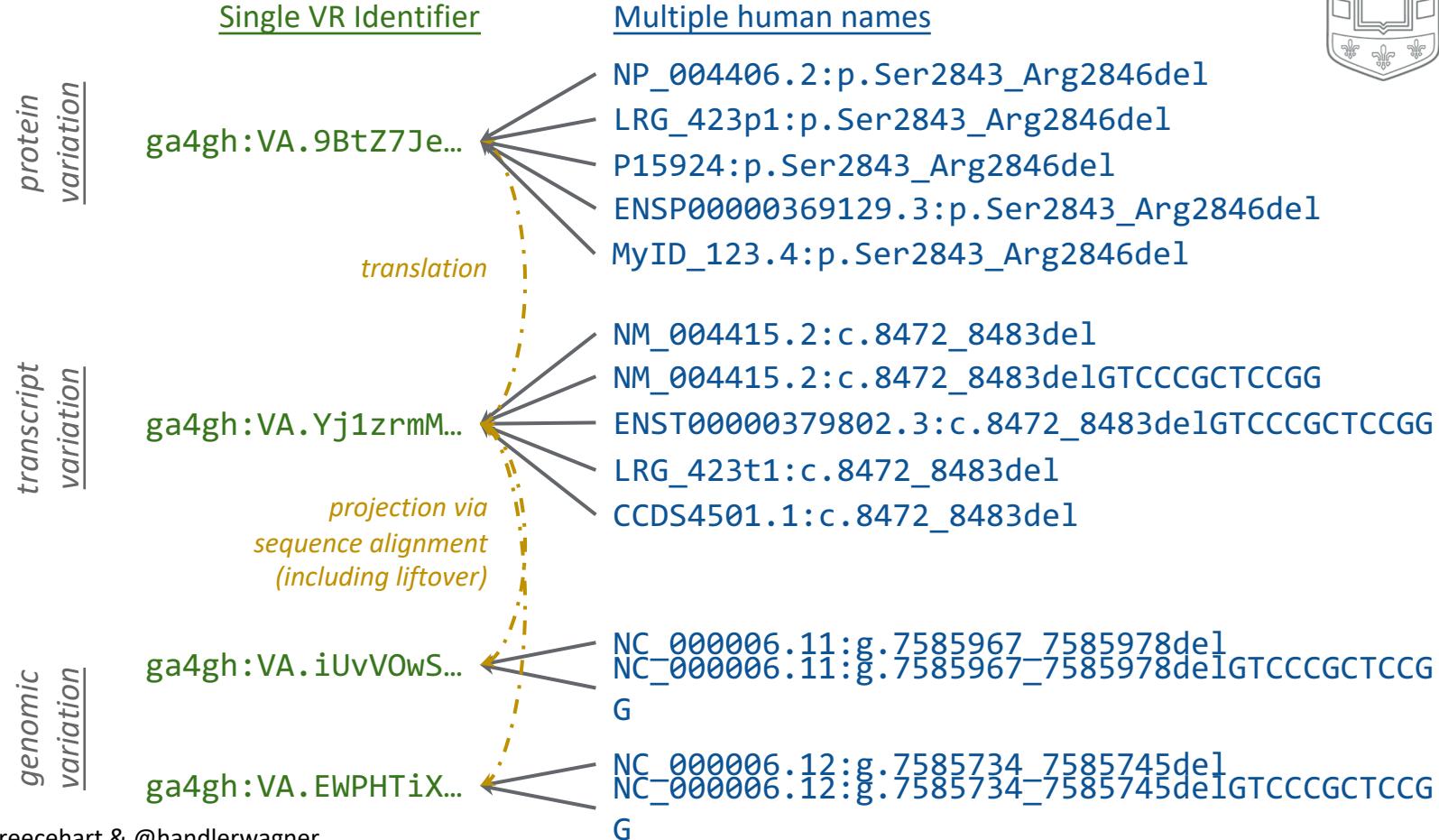
Structured representations
Clear conventions
= VR

One conceptual variant

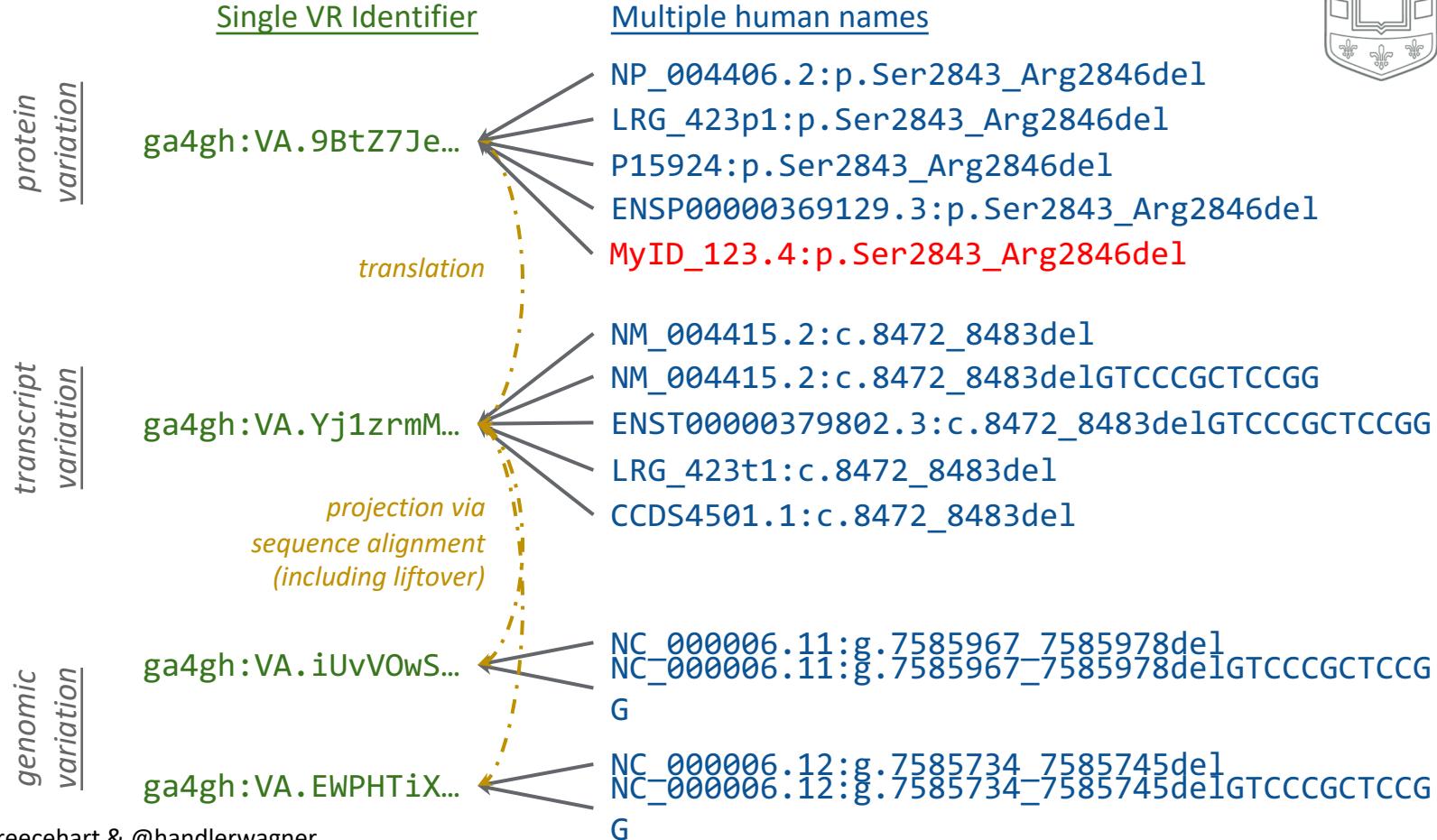
The mission of the VR is to
*standardize the representation of all classes
of biological variation*
to enable
*accurate curation and
reliable computed interpretation.*



VR identifiers: precise, unique names for variation



VR identifiers: precise, unique names for variation



Introduction

Terminology & Information Model

Schema

Implementation Guide

GA4GH Variation Representation Specification

The Variation Representation Specification (VR-Spec) is a standard developed by the Global Alliance for Genomic Health to facilitate and improve sharing of genetic information. The Specification consists of a JSON Schema for representing many classes of genetic variation, conventions to maximize the utility of the schema, and a Python implementation that promotes adoption of the standard.

- [Introduction](#)
- [Terminology & Information Model](#)

<https://vr-spec.readthedocs.io>

Part III – How to Learn Bioinformatics

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

The theme of **learning by doing** is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

The theme of **learning by doing** is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

Learning by doing is not ideal but that's the reality.

-[Mario Caccamo](#)

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

The theme of **learning by doing** is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

Learning by doing is not ideal but that's the reality.

-[Mario Caccamo](#)

The answer depends on **what you want to do**.

-[Manoj Samanta, Homolog.us blog](#)

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

The theme of **learning by doing** is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

Learning by doing is not ideal but that's the reality.

-[Mario Caccamo](#)

The answer depends on **what you want to do**.

-[Manoj Samanta, Homolog.us blog](#)

...get the raw data (SRA) from an a lab evo paper (few muts) and **reproduce results**.

-[Tami Lieberman](#)

Learning bioinformatics **usually requires solving computational problems** of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

The theme of **learning by doing** is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

The answer depends on **what you want to do**.

-[Manoj Samanta, Homolog.us blog](#)

First thing I tell them is **close MS Excel**.

-[Aylwyn Scally](#)

...get the raw data (SRA) from an a lab evo paper (few muts) and **reproduce results**.

-[Tami Lieberman](#)



Learning bioinformatics usually requires solving computational problems of varying difficulty that are extracted from real challenges of molecular biology.

-[Rosalind Bioinformatics Education Platform](#)

Learning bioinformatics is the process of trying something new, making mistakes, and remembering how not to repeat them.

The answer depends on what you want to do.

-[Manoj Samanta, Homolog.us blog](#)

First thing I tell them is close MS Excel.

-[Aylwyn Scally](#)

The theme of learning by doing is probably the one that I suggest most to people.

-[Nicholas J. Loman](#)

Learning by doing is not ideal but that's the reality.

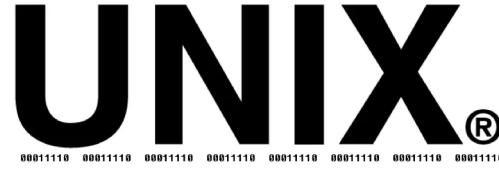
-[Mario Caccamo](#)

...get the raw data (SRA) from an a lab evo paper (few muts) and reproduce results.

-[Tami Lieberman](#)

Trying Something New

Programming



Programming - UNIX

1. Familiarity with OS operations
 - a. creating, moving, deleting files and directories
 - b. I/O streams
 - c. managing permissions
2. Learning common tools
 - a. streaming text manipulation
 - i. sed / awk / tr
 - ii. head / tail
 - iii. more / less
 - b. compression
 - i. tar / gzip
 - ii. tabix
 - c. text editing
 - i. nano
 - ii. emacs / vim
 - d. remote operations
 - i. ssh
 - ii. scp
 - iii. wget / curl



Programming - UNIX

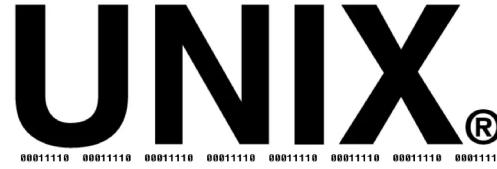
[rnaseq.wiki](#): a broad suite of tutorials covering an end-to-end bioinformatics workflow for RNA-seq analysis.
Includes a UNIX “bootcamp” tutorial:

https://github.com/griffithlab/rnaseq_tutorial/wiki/Unix-Bootcamp

[tutorialspoint](#): generic programming tutorials repository.
Contains a very thorough UNIX introduction module:
www.tutorialspoint.com/unix/index.htm

[Learnshell.org](#): Interactive UNIX shell programming tutorial.

Programming



Python

Programming langua...

Perl

High-level programmi...

R

Programming langua...

+ Add comparison

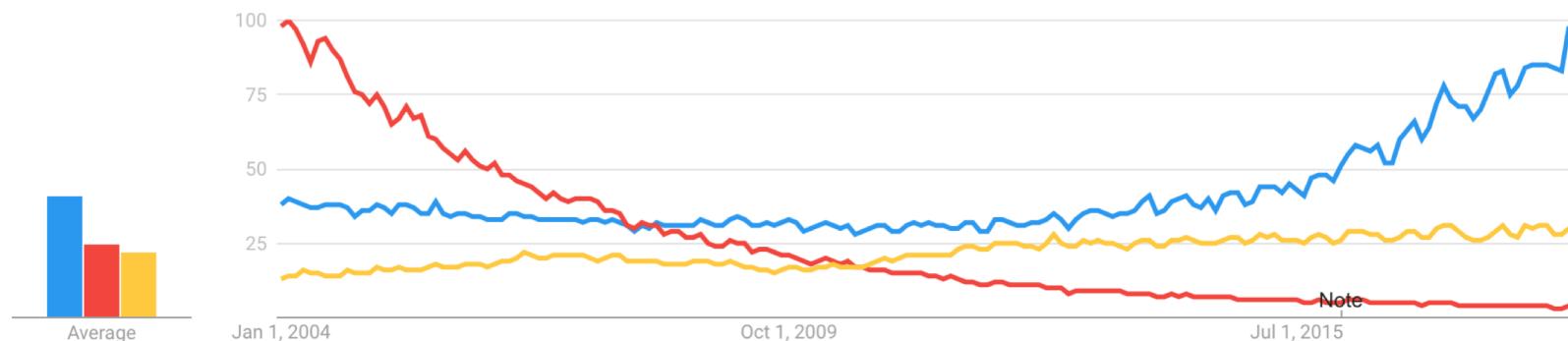
Worldwide ▾

2004 - present ▾

All categories ▾

Web Search ▾

Interest over time



Average

Programming - Python

- General purpose language
 - Readability focused
 - Forgiving interpreter
- Broad applications
 - Machine Learning
 - Scripting
 - Data visualization
 - Web applications
- Great for Reproducibility
 - Environment management
 - Jupyter notebooks
 - PyCharm community IDE



Programming - Python

[Hitchhiker's Guide to Python](#): an opinionated guide to getting started with Python. Includes commonly used extensions.

[ROSLIN](#): bioinformatics learning platform with numerous “learn by doing” exercises, starts with a Python primer.

[The Python Tutorial](#): Provided by the Python foundation, extremely comprehensive survey of the core library.

Programming - Python

Problem

A **string** is simply an ordered collection of symbols selected from some **alphabet** and formed into a word; the **length** of a string is the number of symbols that it contains.

An example of a length 21 **DNA string** (whose alphabet contains the symbols 'A', 'C', 'G', and 'T') is "ATGCTTCAGAAAGGTCTTACG."

Given: A DNA string s of length at most 1000 nt.

Return: Four integers (separated by spaces) counting the respective number of times that the symbols 'A', 'C', 'G', and 'T' occur in s .

Sample Dataset

```
AGCTTTCTTGACTGCAACGGCAATATGTCTCTGTGGATTAAAAAAAGAGTGTCTGATAGCAGC
```

Sample Output

```
20 12 17 21
```

<http://rosalind.info/problems/dna/>

Programming - Python

```
[>>> from collections import Counter
[>>> dna = 'AGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAAGAGTGTCTGATAGCAGC'
[>>> nt_count = Counter(dna)
[>>> nt_count
Counter({'T': 21, 'A': 20, 'G': 17, 'C': 12})
[>>> out_list = [str(nt_count[nucleotide]) for nucleotide in ('A','C','G','T')]
[>>> out_list
['20', '12', '17', '21']
[>>> print(' '.join(out_list))
20 12 17 21
```

<http://rosalind.info/problems/dna/>

Programming - R

- Statistical Computing Language
- Strengths
 - Machine Learning
 - Data visualization
 - ggplot2
- Broad bioinformatics community
 - historical depth with Bioconductor
- Actively maintained IDE
 - RStudio



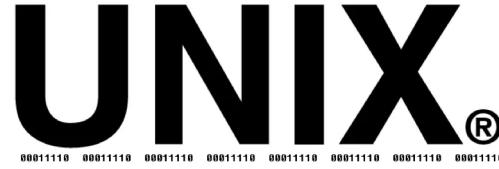
Programming - R

[SWIRL](#): interactive instruction in R programming and data science in the R console.

[R-bloggers](#): R news and tutorials contributed by R blogger community.

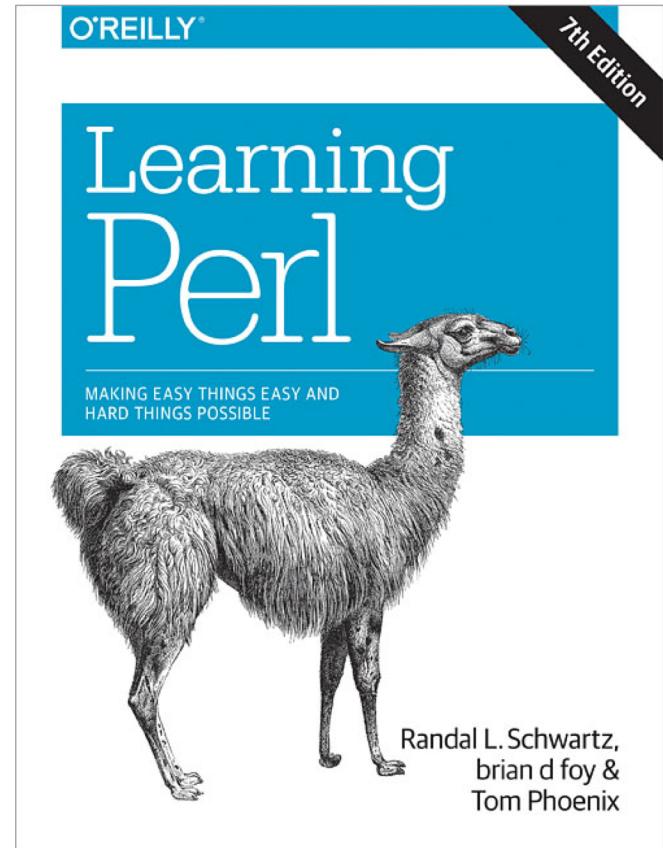
[GenViz.org](#): Genomic Visualizations with R tutorial, contains an R primer.

Programming



Programming - Perl

- Historically popular
- Very useful for text manipulation
- Utility as scripting language in decline



Making Mistakes

Making Mistakes - Resources



Question: Snp, Dip, Snv Notation



A nomenclature question. Considering:

4



- a **SNP** is a **single nucleotide polymorphism**, so a nucleotide change **within-species** of a single position, and
- a **DIP** is a **deletion/insertion polymorphism**, so an **indel within-species** of a single genomic position, how would one call:
- a nucleotide change **between-species** of a single homologous genomic coordinate?
- how about an **indel between-species** of a single homologous genomic coordinate? Is there an official or recognized nomenclature for them?

snp

• 21k views

[ADD COMMENT](#)

•

[link](#)

Not following ▾

modified 3.6 years ago by [scott.a.fay](#) • 0 • written 7.1 years ago by [2184687-1231-83-](#) • 4.9k



do you have a reference to that definition of SNV?



2 [ADD REPLY](#)

• [link](#)

written 7.1 years ago by [Jeremy Leipzig](#) ♦ 17k



I thought the difference between SNP and SNV was one of population frequency (SNP being a frequently observed SNV) not of intra versus inter species differences



1

[ADD REPLY](#)

• [link](#)

written 7.1 years ago by [Russh](#) • 1.2k



It might clarify the SNP and SNV by reading [this blog](#) - "SNP vs. SNP"?



Quote:

11

SNP (single nucleotide polymorphism) vs. SNV (single nucleotide variant) As their name suggests, both are concerned with aberrations at a single nucleotide. However, a SNP is when an aberration is expected at the position for any member in the species – for example, a well characterized allele. A SNV on the other hand is when there is a variation at a position that hasn't been well characterized – for example, when it is only seen in one individual. It is really all a question of frequency of occurrence.



7.1 years ago by

[2184687-1231-83-](#) •

4.9k



7.1 years ago by

[Boboppie](#) • 520

Cambridge, UK

Remembering

Remembering - git



- Speed
- Simple Design
- Strong Support for non-linear development (thousands of branches)
- Fully distributed
- Scalable to very large projects (like the Linux Kernel)

<https://git-scm.com/book/en/v1/Getting-Started>

Remembering - GitHub



- Hosts Git Repositories
- Documentation Markdown
- Wikis
- Issue Tracking
- Repo Management
- Collaborative Workflow



ahwagner / **civic_pmids.ipynb**

Created 10 months ago

Edit

Delete

Star 0

Code

Revisions 1

Embed ▾

<script src="https://gist.



Download ZIP

Jupyter Notebook for extracting all PMIDs in ClViC

civic_pmids.ipynb



Raw

In [1]: import requests

In [2]: url = "https://civic.genome.wustl.edu/api/sources"

```
In [3]: pmids = set()
next_page = url
while next_page:
    resp = requests.get(next_page)
    resp.raise_for_status()
    json = resp.json()
    for record in json['records']:
        pmid = record.get('pubmed_id', None)
        if pmid is None:
            continue
        pmids.add(pmid)
    next_page = json['meta']['links']['next']
```

In [4]: len(pmids)

Out[4]: 1752

```
In [5]: with open('pmids.txt', 'w') as out:
    for pmid in pmids:
        print(pmid, file=out)
```

Remembering - GitHub



- Hosts Git Repositories
- Documentation Markdown
- Wikis
- Issue Tracking
- Repo Management
- Collaborative Workflow

<https://guides.github.com/activities/hello-world/>



Learning bioinformatics usually requires solving computational problems of varying difficulty that are extracted from real challenges of molecular biology.

-Rosalind Bioinformatics Education Platform

Learning bioinformatics is the process of trying something new, making mistakes, and remembering how not to repeat them.

The answer depends on what you want to do.

-Manoj Samanta, Homolog.us blog

First thing I tell them is close MS Excel.

-Aylwyn Scally

The theme of learning by doing is probably the one that I suggest most to people.

-Nicholas J. Loman

Learning by doing is not ideal but that's the reality.

-Mario Caccamo

...get the raw data (SRA) from an a lab evo paper (few muts) and reproduce results.

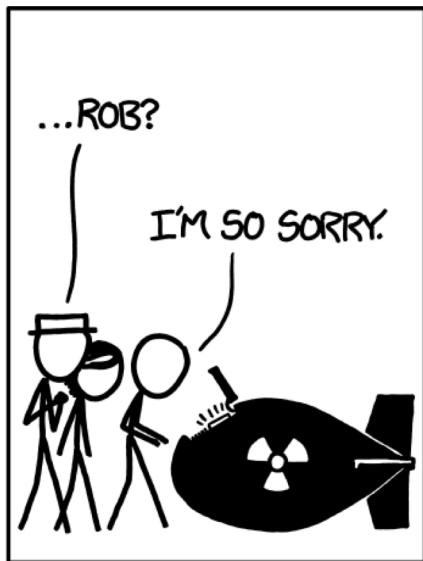
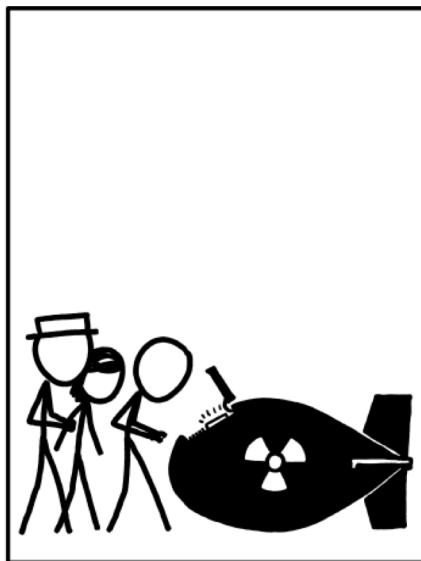
-Tami Lieberman



First thing I tell them is **close MS Excel**.

-[Aylwyn Scally](#)

Questions?



<https://xkcd.com/1168/>