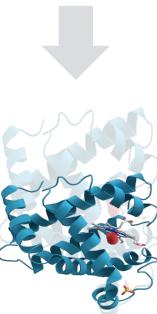


AlphaMissense predictions for human genetic variation research

gtatatatcatgatatgtatacgcatacatgagctca



Webinar objectives

- *Describe the principles behind AlphaMissense and its approach to variant classification*
- *Identify how AlphaMissense utilises protein structure and sequence information for variant assessment*

Talk outline

- **Overview of AlphaMissense**
- Design principles
- Evaluating predictions
- Using the scores

Note: Key slides designated with this:



Google DeepMind and science

- Mission: **build AI responsibly to benefit humanity**
- Science group which approaches fundamental science problems
- One such problem is genome variant Interpretation
- Definition: A **missense variant** is a variant affecting a protein, by changing a reference amino acid to an alternate amino acid.

Reference:

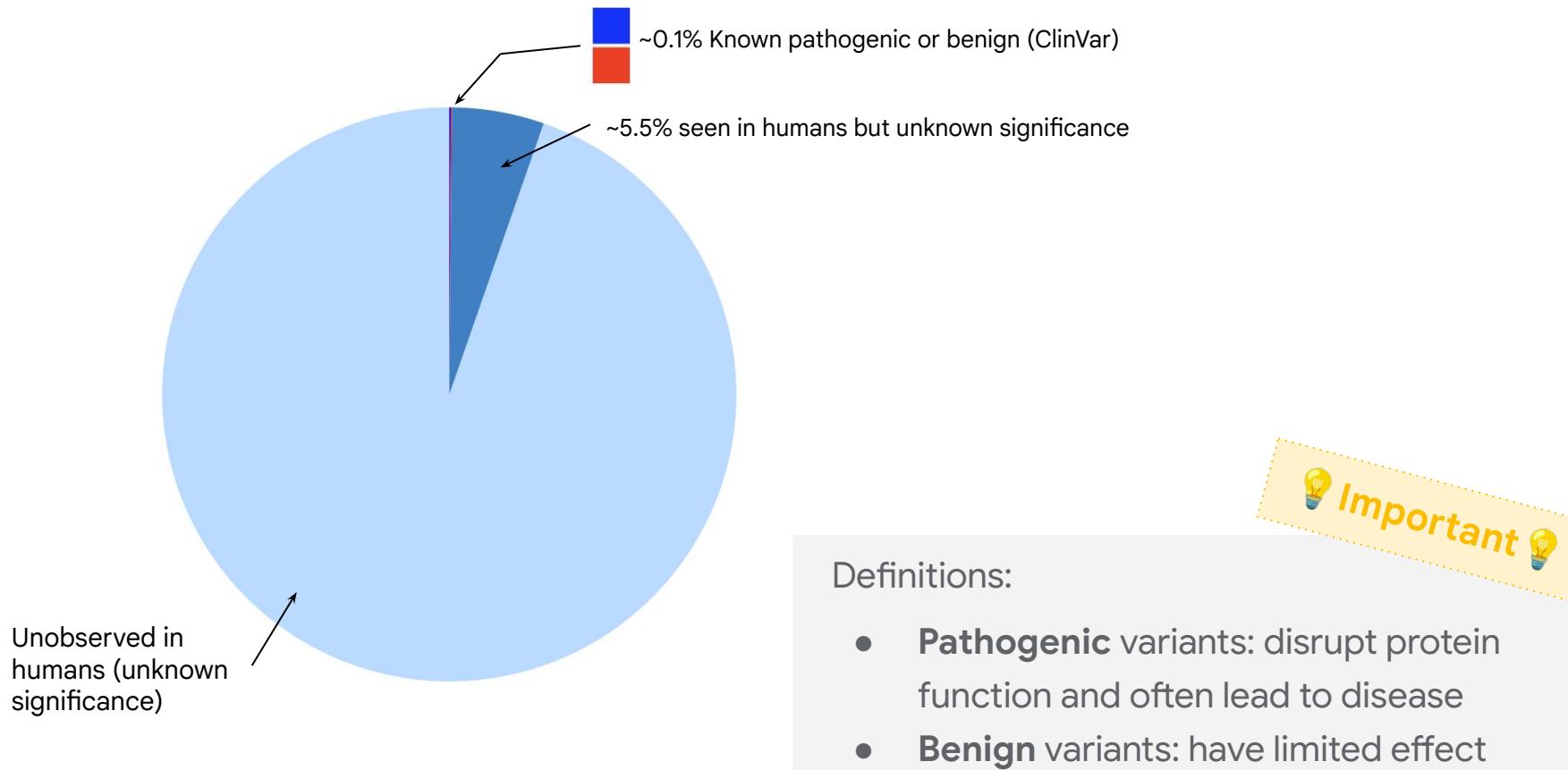
CAG ----- DNA
MDVVAMVNQTVATMIS ----- Protein



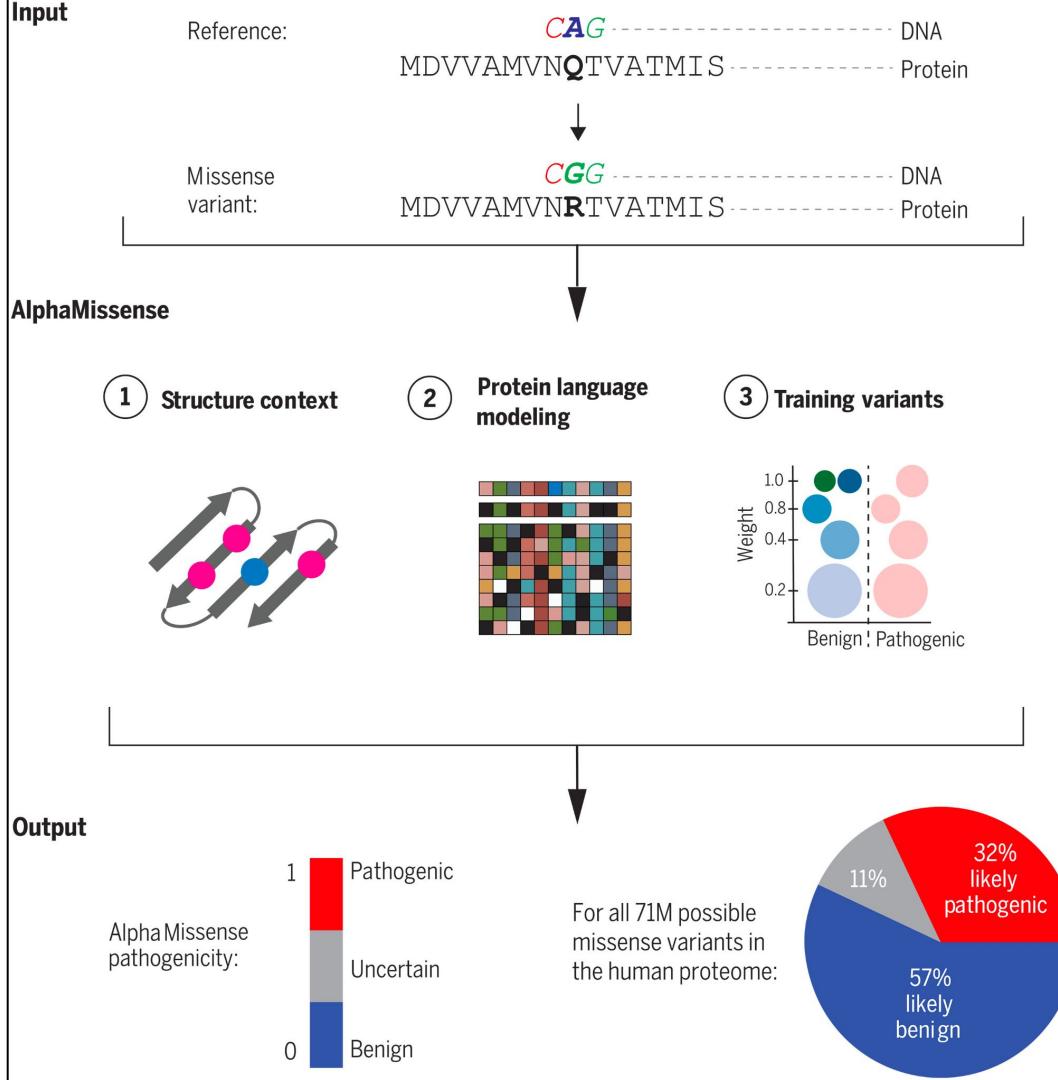
Missense
variant:

CGG ----- DNA
MDVVAMVNRTVATMIS ----- Protein

Why missense? A major gap in our knowledge of variant effects



AlphaMissense: proteome-wide missense variant effect prediction



Webinar objectives

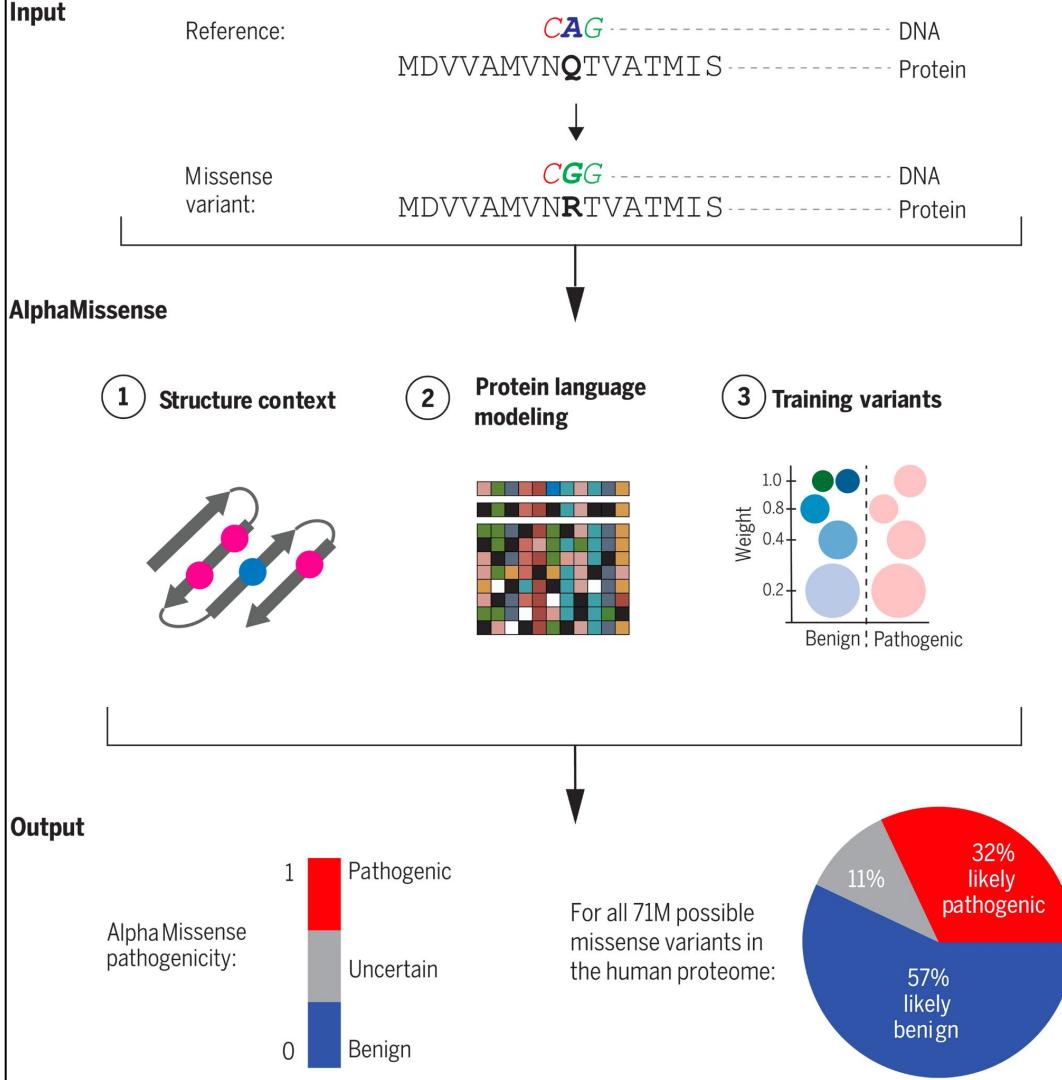
- *Describe the principles behind AlphaMissense and its approach to variant classification*
- *Identify how AlphaMissense utilises protein structure and sequence information for variant assessment*

Talk outline

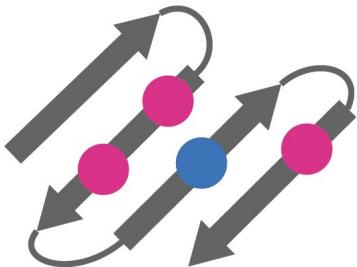
- Overview of AlphaMissense
- **Design principles**
- Evaluating predictions
- Using the scores

Why is machine learning an appropriate choice for this problem?

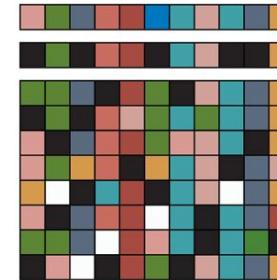
- **Clear task**
 - Input: Variant
 - Output: Score from 0-1
- **Clear metrics**
 - Held out ClinVar variants
 - Deep mutational scans
- **Clear impact**
 - Fill in the gaps in an important database
- **Clear inductive biases**



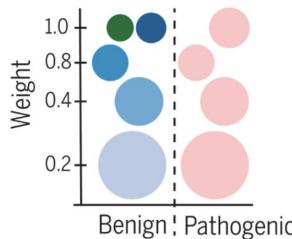
What inductive biases (principles) inform the model?



- Protein context matters
 - Sequence
 - Structure

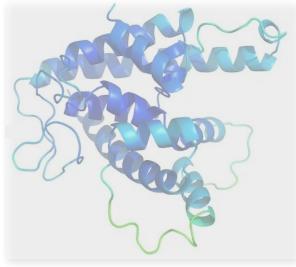


- Protein evolutionary history matters

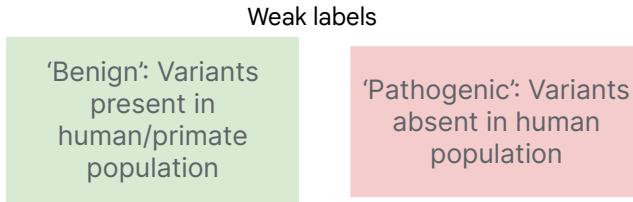


- Observed sequence variation in humans matters

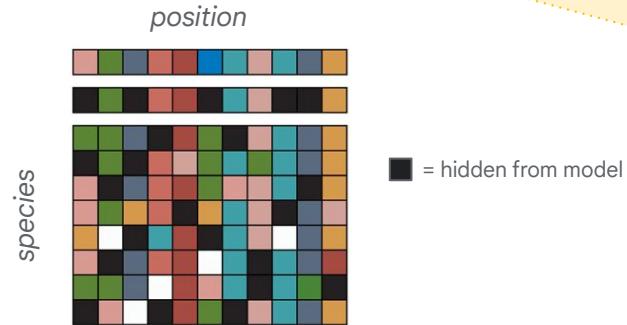
Under the hood: Adapting AlphaFold into AlphaMissense



AlphaFold2 takes protein sequence as input and returns its **structure**. This is the *starting point* for AlphaMissense.



AlphaMissense fine tunes this likelihood via an additional challenge: given a human protein variant, classify whether it is present or absent in the human population, thereby weakly supervising on **human sequence variation**.



AlphaFold2 has a component that predicts the likelihood of amino acids belonging in *hidden entries* in the multiple sequence alignment of the input protein - thereby learning the statistics of **protein evolutionary history**.



Finally, AlphaMissense is evaluated on ClinVar to assess performance. Note: ClinVar variants are EXCLUDED from the previous step. The model is NOT trained on ClinVar data directly.

Webinar objectives

- *Describe the principles behind AlphaMissense and its approach to variant classification*
- *Identify how AlphaMissense utilises protein structure and sequence information for variant assessment*

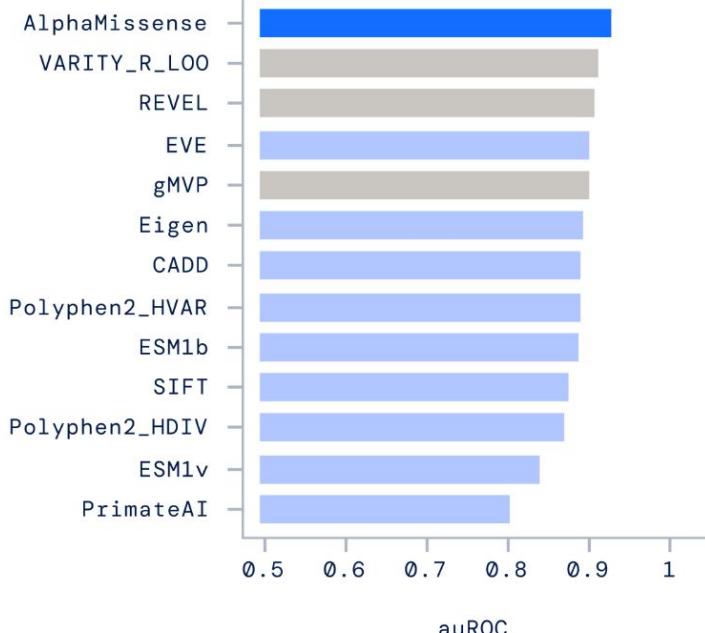
Talk outline

- Overview of AlphaMissense
- Design principles
- **Evaluating predictions**
- Using the scores

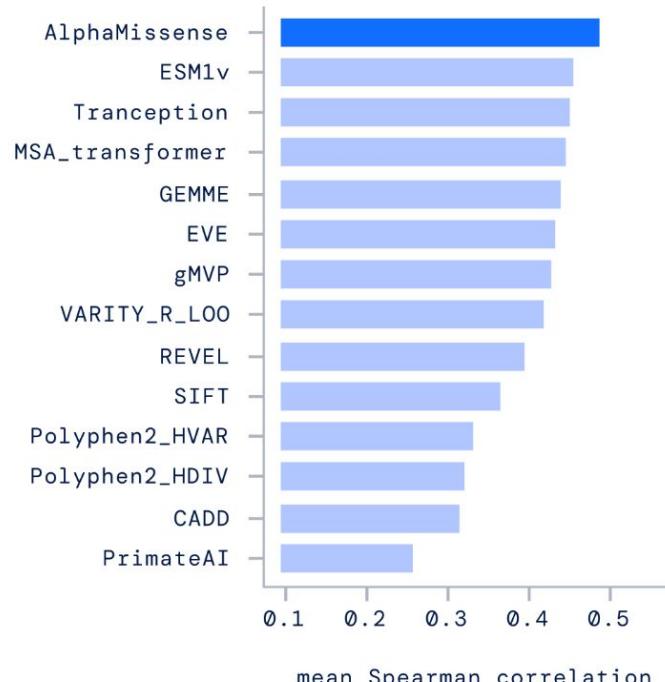
Good performance across a diverse set of benchmarks

ClinVar (Class-balanced 18924 variants)

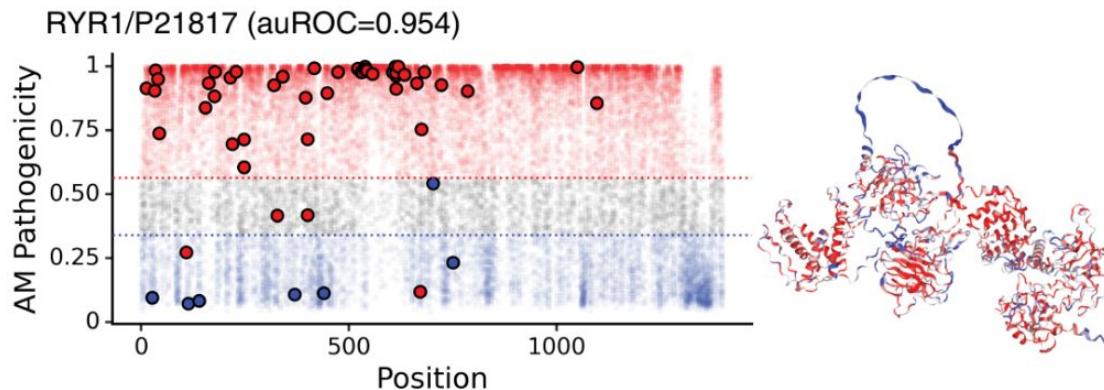
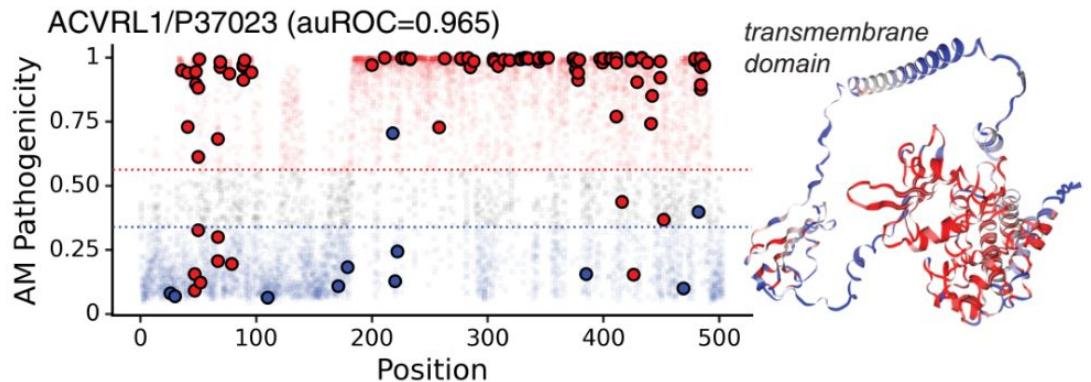
Trained on ClinVar



Experimental assays (25 proteins)



Example: AlphaMissense predictions for clinically actionable genes



All about the 0 to 1 score



What does a score between 0-1 mean?

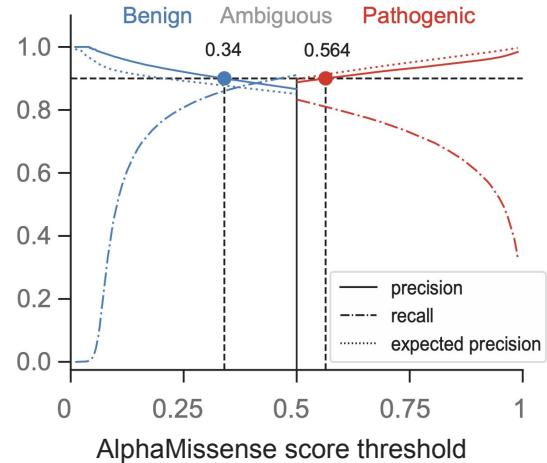
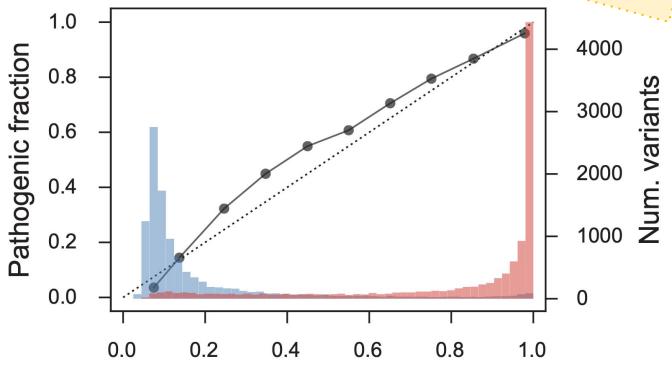
- Score of 0.8 implies: for 10 variants that score in that range, you will expect 8 of them will say “pathogenic/likely-pathogenic” in ClinVar, and 2 of them to say something else (uncertain or likely benign)

How did they get to be in that scale?

- Model output was transformed to [0,1] using a logistic regression classifier on a set of ClinVar validation variants

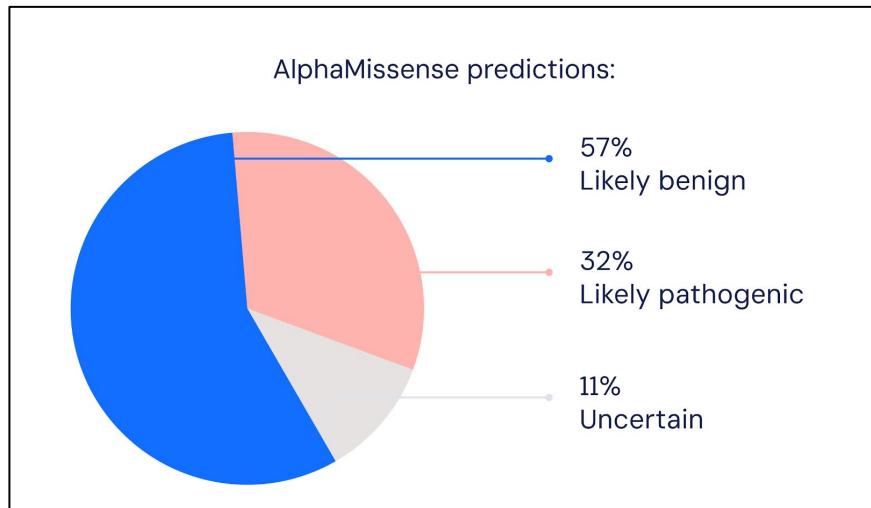
How were class thresholds drawn?

- Likely benign [0-0.34]: The set of predictions in this range obtain 90% negative-class precision (benign)
- Likely pathogenic [0.564-1]: Same, but for 90% positive-class precision (pathogenic).
- Class labels chosen to match ClinVar naming conventions - however, these should not be interpreted as equivalent to ACMG labels



Proteome-wide predictions with more confidently classified variants

- **600M predictions** all AA substitutions and missense variants in 20k canonical gene isoforms and 60k alternative isoforms
- **Higher coverage %** of confident predictions due to better performance. Fractions are based on thresholds that yield 90% precision on ClinVar variants



Webinar objectives

- *Describe the principles behind AlphaMissense and its approach to variant classification*
- *Identify how AlphaMissense utilises protein structure and sequence information for variant assessment*

Talk outline

- Overview of AlphaMissense
- Design principles
- Evaluating predictions
- **Using the scores**

How should I reason about the scores?



If you are a clinical researcher:

- Scores for genetic variants are available via familiar portals such as VEP, DECIPHER
- There is a lot of thinking and guidelines around use of informatic evidence in clinical settings that should inform use of AlphaMissense (and other predictive approaches)
 - AlphaMissense scores are a research tool
 - Should be used alongside other sources of evidence when drawing scientific conclusions
 - Not a substitute for professional medical advice, diagnosis or treatment
- Recommend benchmarking for your specific research use-case



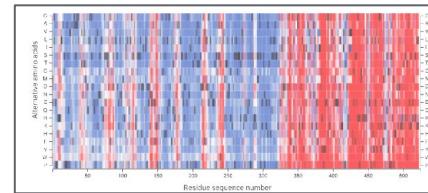
How should I reason about the scores?



If you are a molecular biologist:

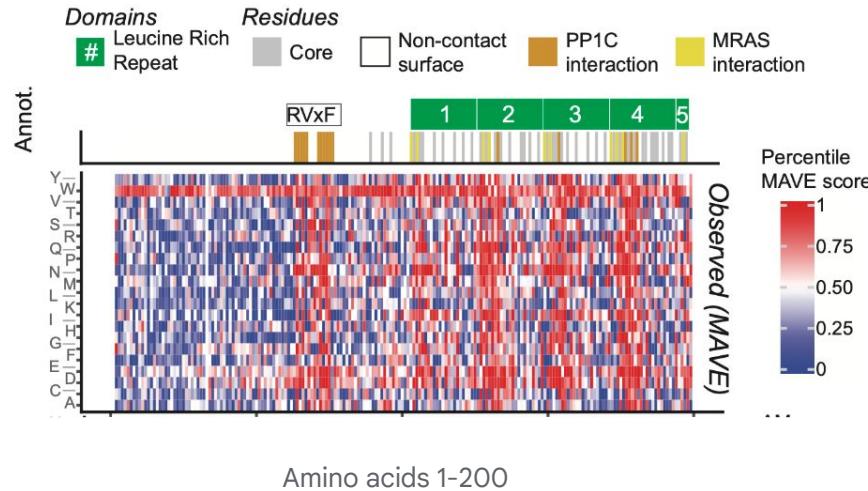
Recommend three levels of granularity with which to approach the scores:

- Individual protein - all variants
 - Residue position x 20 possible amino acids (1 ref, 19 alt)
 - Heatmap visualization recommended
- Individual protein - position averages
 - Single value for each position in your protein
 - Use in conjunction with protein structure (measured or predicted) to aid in interpretation
- Protein-level averages



Is AlphaMissense informative of function? Example SHOC2

Experiment (cancer cell fitness) Kwon et al, 2022, Nature

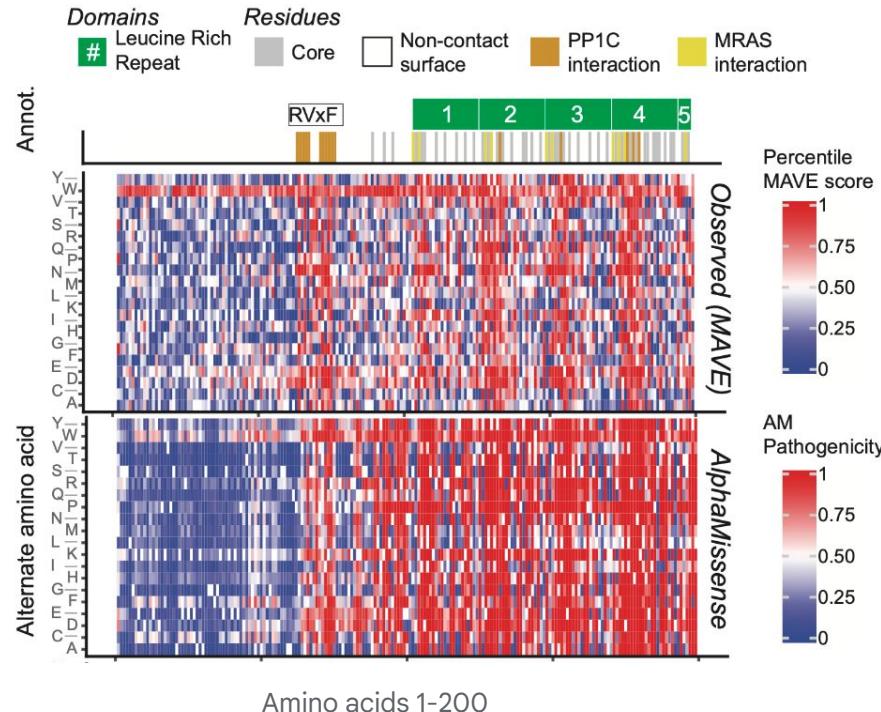


Is AlphaMissense informative of function? Example SHOC2

Experiment (cancer cell fitness) Kwon et al, 2022, Nature

AlphaMissense

Spearman R = 0.47



Is AlphaMissense informative of function? Example SHOC2

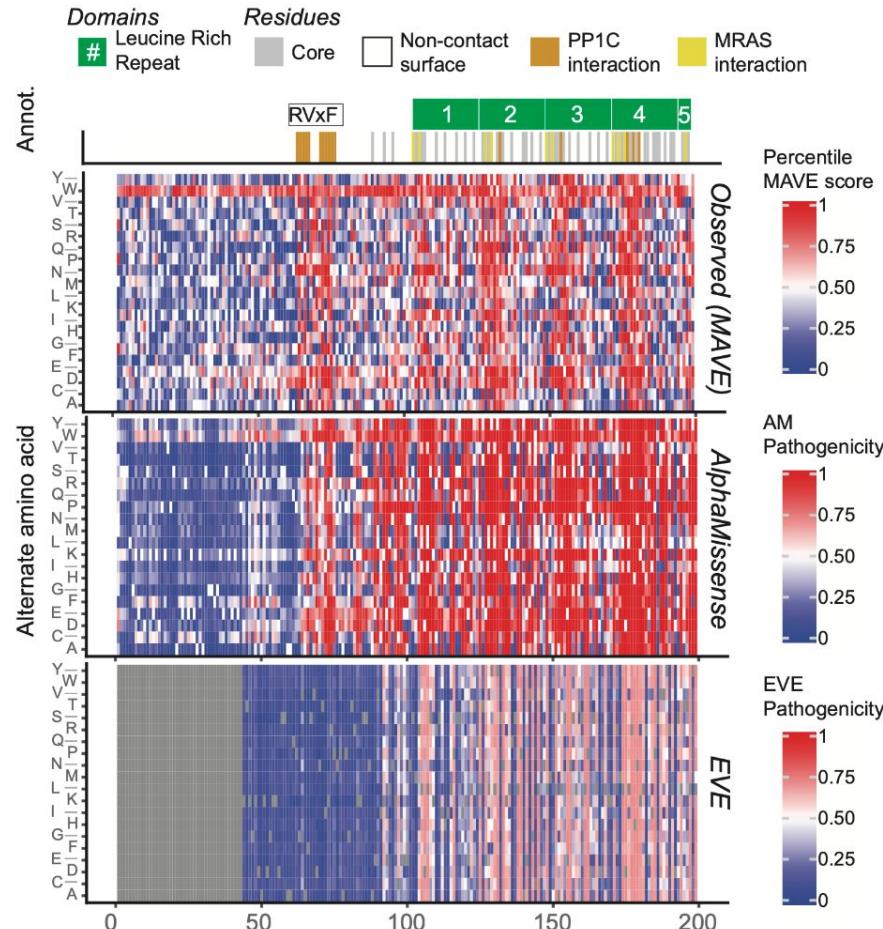
Experiment (cancer cell fitness) Kwon et al, 2022, Nature

AlphaMissense

Spearman R = 0.47

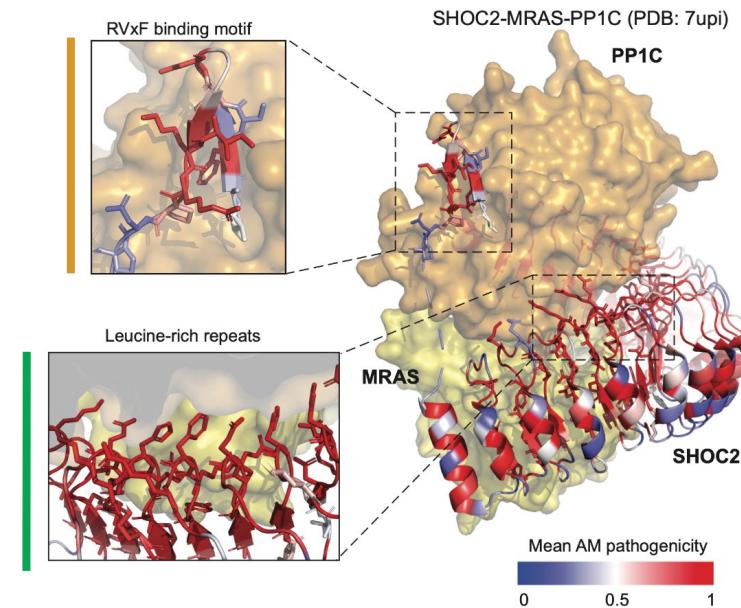
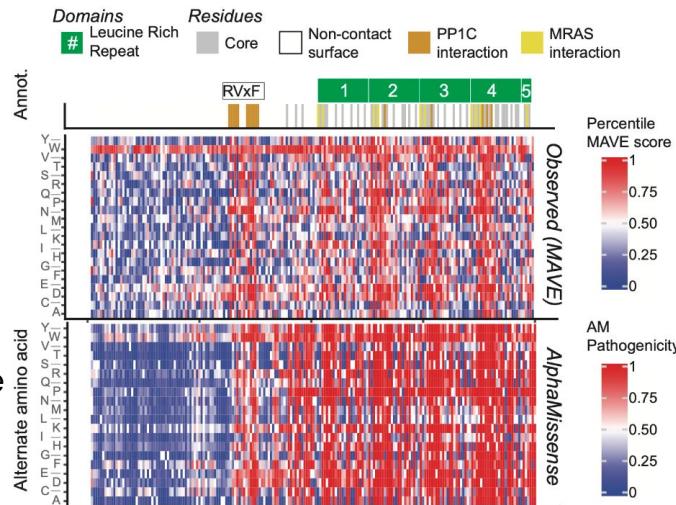
EVE

Spearman R = 0.32



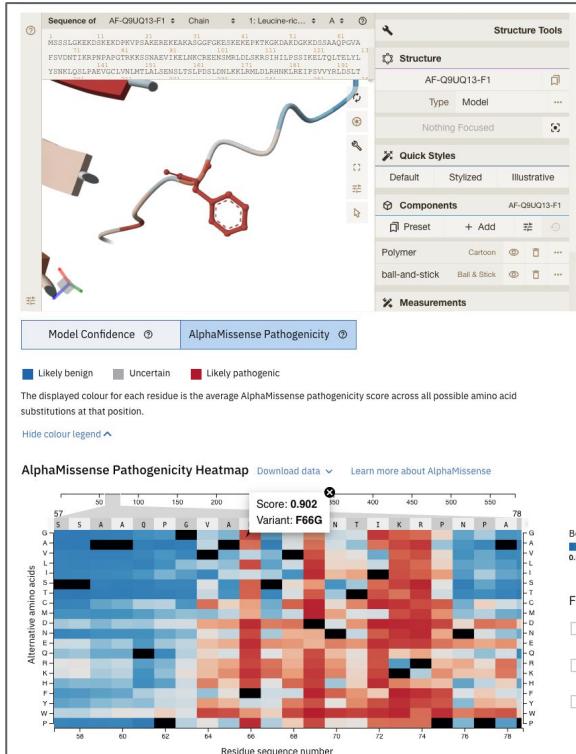
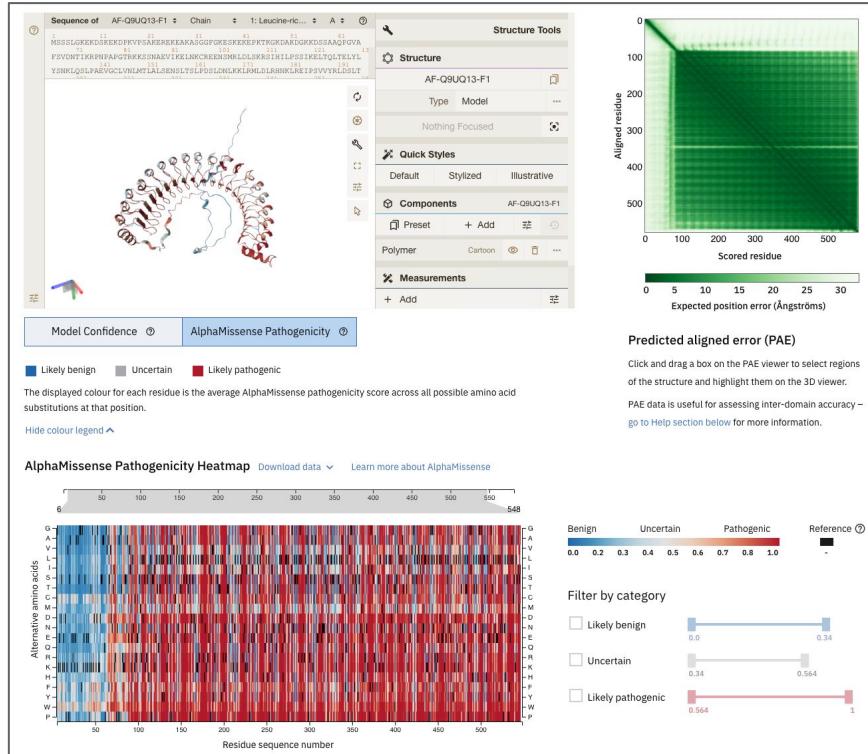
Is AlphaMissense informative of function? Example SHOC2

Experiment
(cancer cell
fitness)



- Regions with high AlphaMissense predictions reflect recently discovered functionally-significant domains.

Try it yourself using AlphaFold Database!



<https://alphafold.ebi.ac.uk/entry/Q9UQ13>

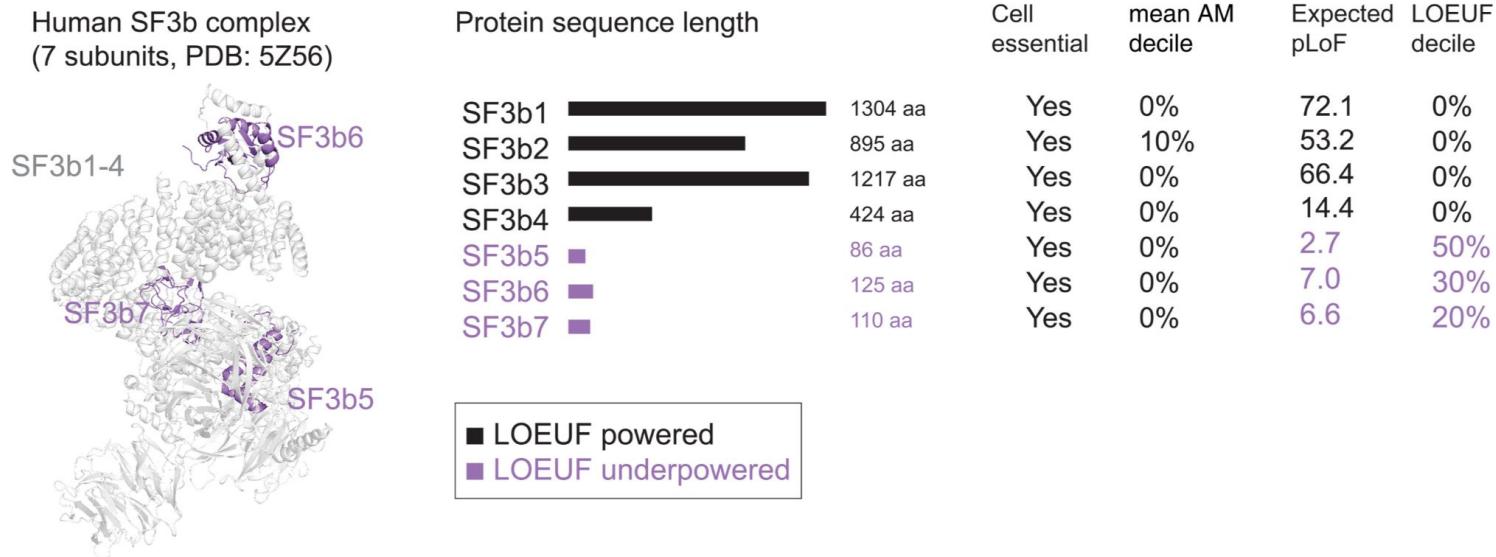
How should I reason about the scores?

If you are a molecular biologist:

There are three levels of granularity with which to approach the scores:

- Individual protein - all variants
 - Residue position x 20 possible amino acids (1 ref, 19 alt)
 - Heatmap visualization recommended
- Individual protein - position averages
 - Single value for each position in your protein
 - Use in conjunction with protein structure (measured or predicted) to aid in interpretation
- **Protein-level averages**
 - Single number for your protein of interest
 - Good for estimating the protein intolerance to loss-of-function variation, LOEUF-style

Predicting gene essentiality by averaging AM scores



Limitations!



Shares some of the same limitations of AlphaFold2

- Certain length restrictions
- Certain non-canonical residues
- Lack of context - PPI, DNA, RNA, etc.

The 0 to 1 score does not inform on the following:

Clinical:

- Penetrance
- Haploinsufficiency
- Specific diseases or phenotypes
- Epistasis (combinations of variants)

Molecular biology:

- Biochemical details (stability vs. function)
- Not directly equivalent to a mutational experiment (which are tied to specific biochemical readout), though the information can be complementary
- Does not return the structures of proteins containing the alternative amino acid

Other notes:

- Gain of function variants are challenging
- Thresholds might need to be tweaked for particular use cases

Not a software: main artifact is the database

Webinar objectives

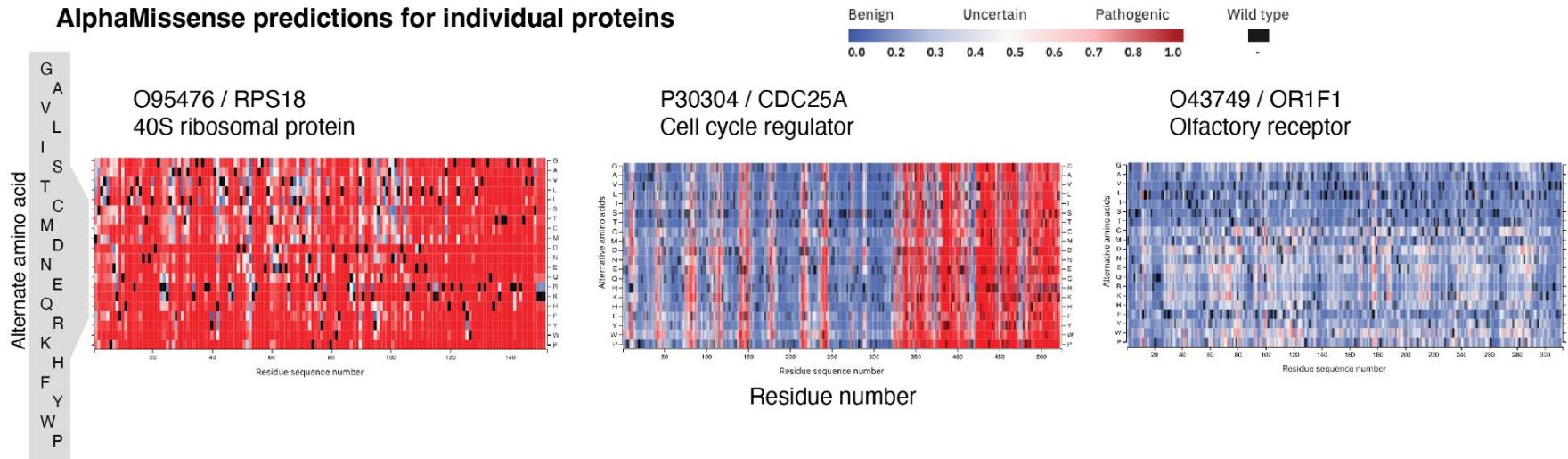
- *Describe the principles behind AlphaMissense and its approach to variant classification*
- *Identify how AlphaMissense utilises protein structure and sequence information for variant assessment*

Talk outline

- **Bonus: Field notes**

Field notes

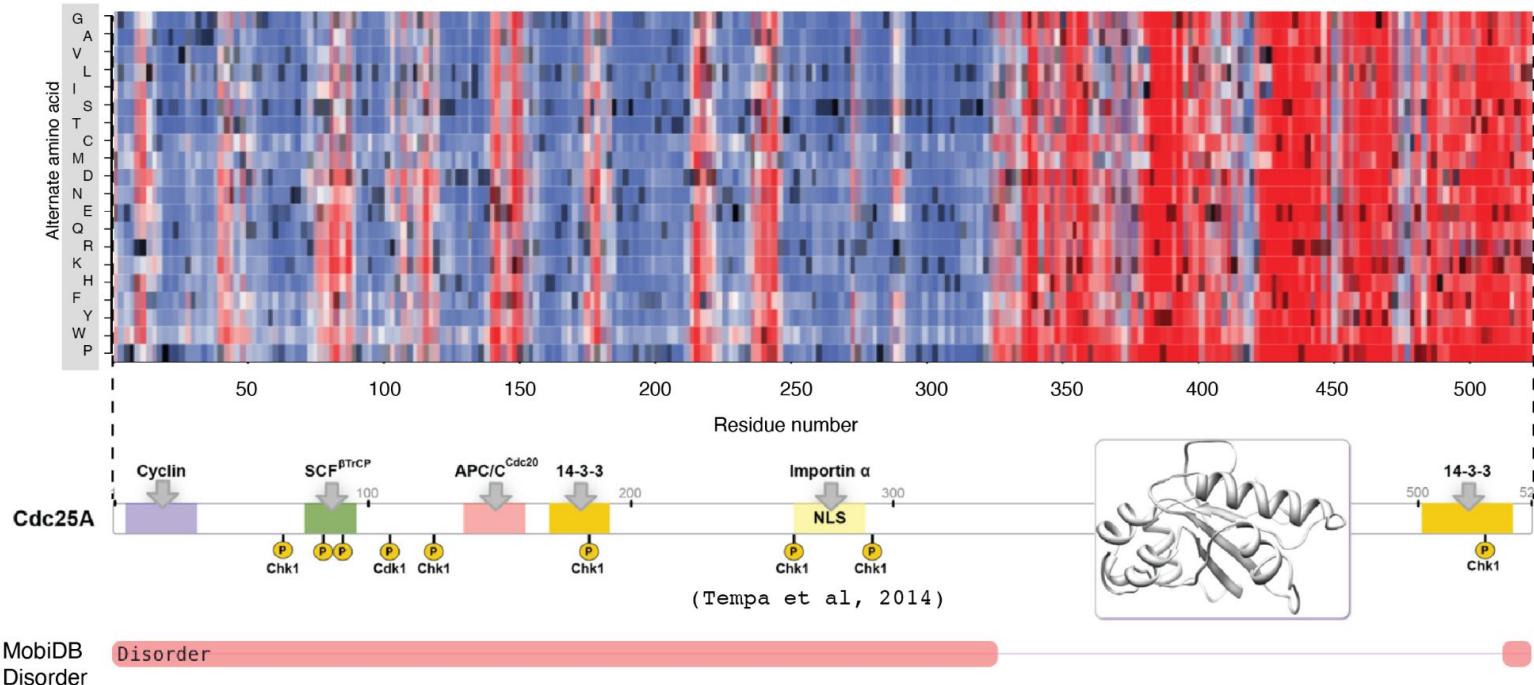
(Note: These are anecdotes, not from the peer reviewed paper.)



Field notes

(Note: These are anecdotes, not from the peer reviewed paper.)

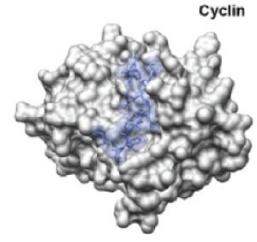
Closeup for CDC25A



Field notes - interaction motifs

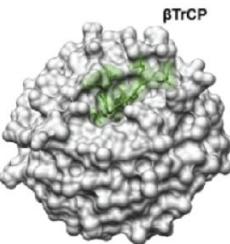
(Note: These are anecdotes, not from the peer reviewed paper.)

Individual short linear interacting motifs



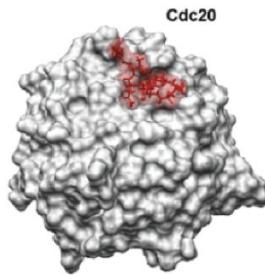
Cyclin
CDC25A 1st H R R R L L F A C 1st
CDKN1B 1st S A C R N L F G P 1st
CDC6 1st V K G R R L V F D N 1st
E2F1 1st V K R R L D E E T 1st
TP53^{1st} M I R H K K L M F K T 1st
RB1 1st L K K L R F D I 1st

Consensus | [KR]xLxP-x[ILVF]



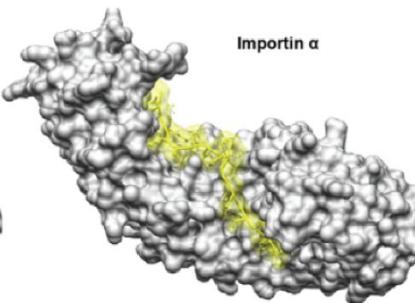
βTrCP
CDC25A 1st S T D S G F C L D S P 1st
CTNNB1^{1st} Y L D S G I H S G A T 1st
ATF4 1st D N D S G I C M S P E 1st
DLG1 38th T K D S G L P S Q G L 102th
FBXO5 1st Y E D S G Y S S F S 1st
NFKBIB 1st W C D S G L G S L G P 1st

DpSGxQ-xpS

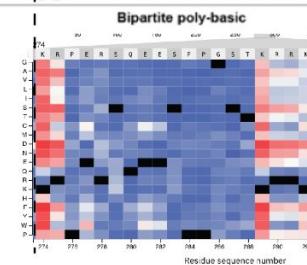


Cdc20
CDC25A 1st E N K E N E A 1st
BUB1B^{1st} I L S K E N Q V I 1st
CDC20 1st L S K E N Q P 100th
CDC 1st G K K E N G P 1st
CDCAS 1st L E K E N E P 1st
PTTG1 1st V D K E N G E 1st

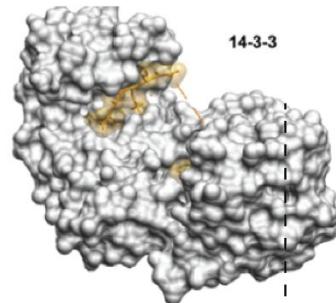
KEN



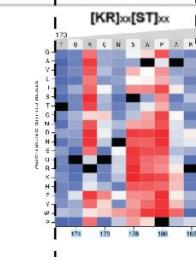
Importin α
CDC25A 1st K R P E R S Q E E S P P G S T K R R K S 1st
RBI^{1st} L K R S A E G S N P P K P L K K L R F D 177th
CDKN1A 1st I K R R Q T S M T D F Y H S K R R L I F S I 258th
N1/N2^{1st} S M I K R K T E E S P L K D K D A K K S K Q I 333th
NPM^{1st} V K R P A A T K K A G Q A K K K K L D K 171th
TP53 1st K R A L P N N T S S S P Q P K K K P L D 223th



(Tempa et al, 2014)



14-3-3
CDC25A 1st I Q R Q N S A P A R 100th
CDC25A 104 R T K S R T W A G E 109th
PRKCE^{1st} E D R S K S A P T S I 148th
PRKCE^{1st} R K A L S F D N R I 170th
ATXN1 372 R K R R W S A P E S 177th
FOXO4 139 R R R A A S M D S S 199th

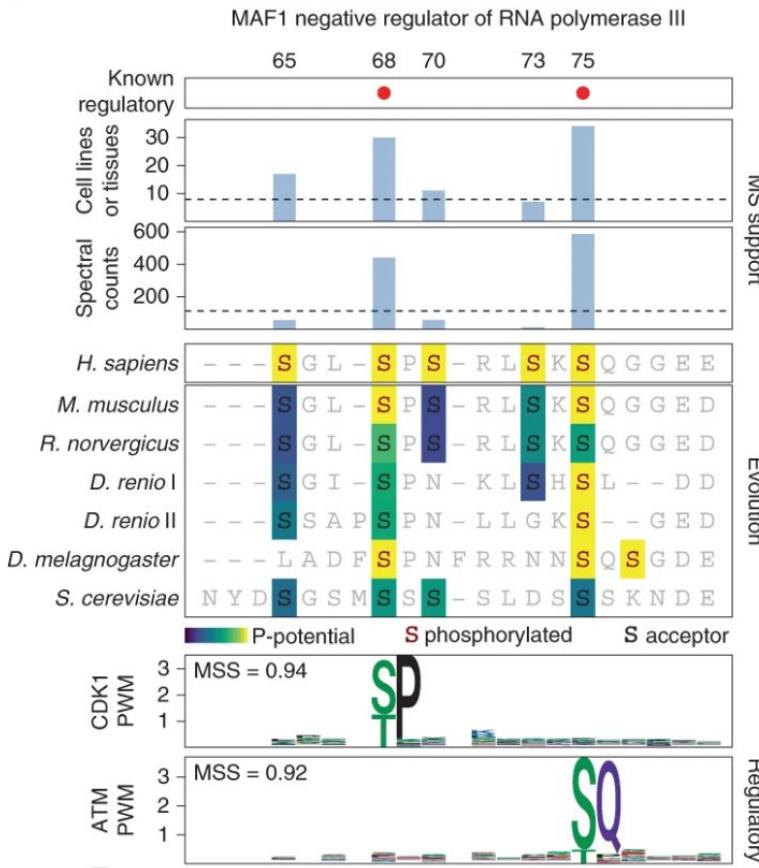
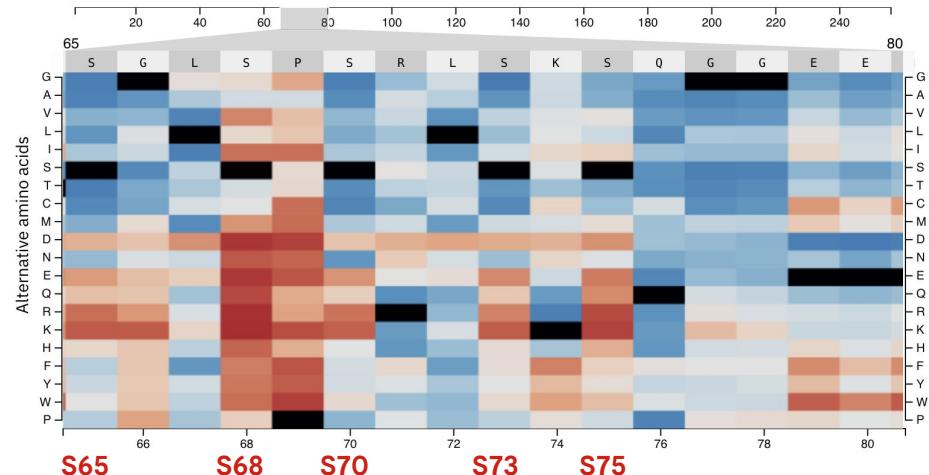


Field notes - phosphorylation sites

Ochoa et al, 2019

(Note: These are anecdotes, not from the peer reviewed paper.)

AlphaMissense Pathogenicity Heatmap [Download data](#) [Learn more about AlphaMissense](#)



Summary

- Fine-tuned AlphaFold to predict pathogenicity of missense variants, without using clinically-ascertained variants in training.
- Outperforms state-of-the-art on multiple diverse benchmarks.
- Increased the number of confidently classified variants (using ClinVar to estimate precision) proteome-wide
- Widely available predictions (Stay on the webinar to learn how to access!)

Citation

Cheng *et al.*, Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Science* **381**, eadg7492(2023). DOI:[10.1126/science.adg7492](https://doi.org/10.1126/science.adg7492)

Acknowledgements

AlphaMissense

- Jun Cheng
- Guido Novati
- Clare Bycroft
- Akvilė Žemgulytė
- Taylor Applebaum
- Alexander Pritzel
- Lai Hong Wong
- Michal Zielinski
- Tobias Sargeant
- Rosalia G. Schneider
- Andrew W. Senior
- John Jumper
- Demis Hassabis
- Pushmeet Kohli
- Žiga Avsec

Contact us:

alphamissense@google.com

Thank you.

Integration of AlphaMissense data into EMBL-EBI resources

7th June 2024

Paulyna Magaña
Bioinformatician



AlphaMissense in Ensembl

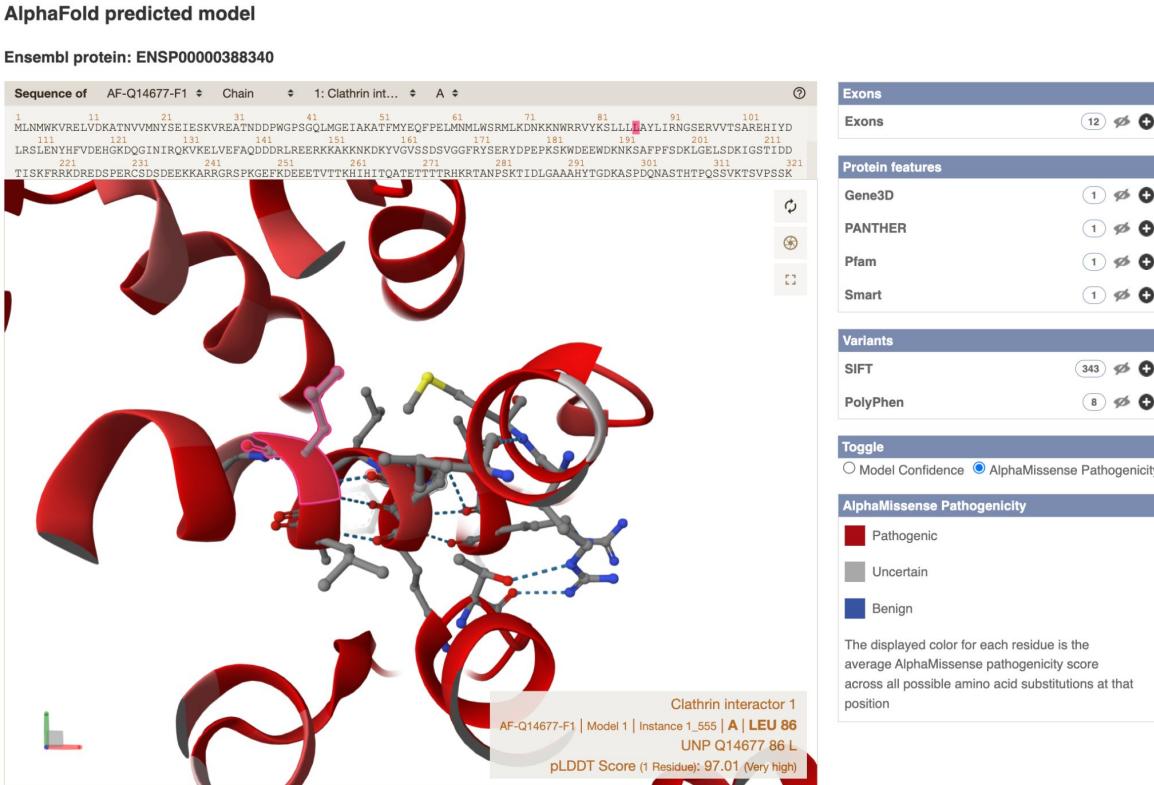
- The Ensembl genome browser (www.ensembl.org) supports research in comparative genomics, evolution, sequence variation and transcriptional regulation.
- AlphaMissense scores are integrated into the Ensembl Variant Effect Predictor (VEP) tool enabling annotation of variants via its web interface, REST API and the command line.



Location	Allele	Consequence	Symbol	Protein position	Amino acids	SIFT	PolyPhen	Pubmed	AlphaMissense classification	AlphaMissense pathogenicity score
7:117587806-	A	missense_variant	CFTR	551	G/D	0.01	0.999	73 PubMed IDs	likely_pathogenic	0.9897
117587806	T	missense_variant	CFTR	551	G/V	0	1	73 PubMed IDs	likely_pathogenic	0.9593
10:62813413-	A	missense_variant	EGR2	409	R/W	0	0.996	20301384, 9537424, 12525712	likely_pathogenic	0.9952
62813413	A	missense_variant	EGR2	409	R/W	0	0.996	20301384, 9537424, 12525712	likely_pathogenic	0.9952
14:73219188-	G	missense_variant	PSEN1	435	L/V	0	-	-	likely_pathogenic	0.9848
73219188	T	missense_variant	PSEN1	435	L/F	0	-	-	likely_pathogenic	0.9978

AlphaMissense in Ensembl

- Average AlphaMissense pathogenicity scores for each amino acid can be visualised on the AlphaFold predicted 3D protein structure, available from the associated Ensembl transcript page.
- The interactive view allows switching between variants, domains, exons and AlphaMissense results to support interpretation of different regions of the protein and their sensitivity to change.



AlphaMissense in DECIPIER

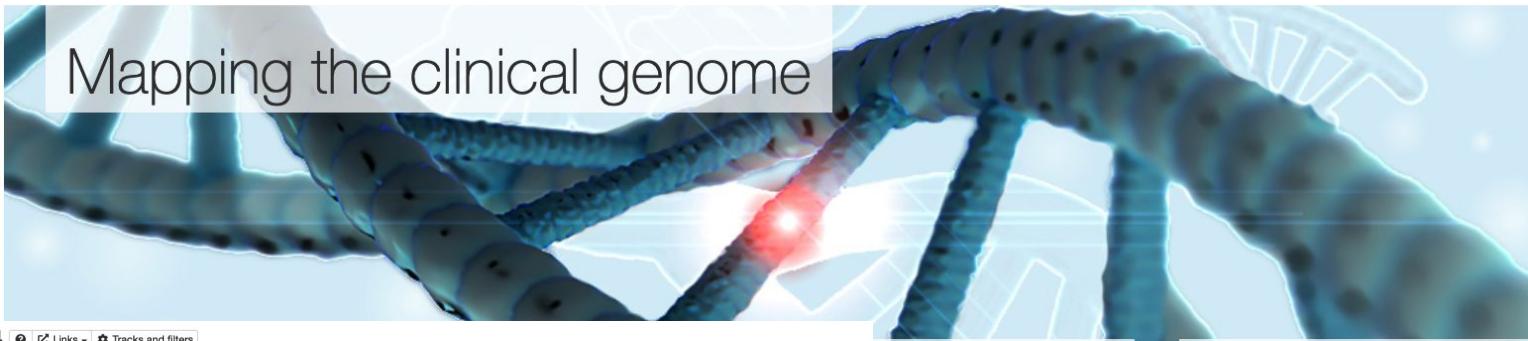


About Browse DDD (UK)

12-6596097-G-A



Mapping the clinical genome



CDK13: Q14004

Links Tracks and filters

Transcript used in protein view: ENST00000181839.10, NM_003718.5 14 exons, 1512aa MANE Select

Exons
Predicted NMD Escape
Conservation
Missense Constraint
gnomAD Missense
gnomAD Homozygous Missense
gnomAD LoF

DECIPHER variants

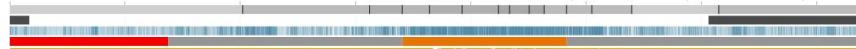
Domains and annotations

ClinVar variants

Secondary Structure

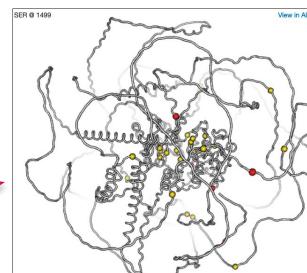
3D Structures

1aa 200aa 400aa 600aa 800aa 1,000aa 1,200aa 1,400aa



DECIPHER enables the sharing and interpretation of phenotype-linked genomic variants.

It aggregates the latest information about variants and makes it available via easy-to-use tools and visualisations.



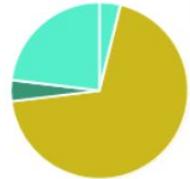
DECIPHER

Gene Browser Protein Annotation Matching patient variants 107 Matching DDD research variants 0

Consequence prediction (VEP) Allele frequency Functional ClinVar Disease cohorts

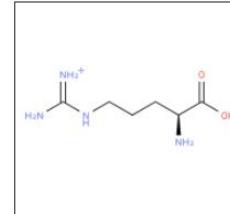
Ensembl Variant Effect Predictor (VEP)

Consequences

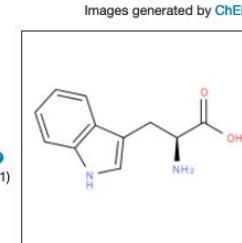


- downstream_gene_variant: 4%
- missense_variant: 69%
- non_coding_transcript_exon_variant: 4%
- upstream_gene_variant: 23%

Amino acid substitution



Grantham distance ?
Moderately radical (101)



Images generated by ChEBI

Annotations for CHD4: 1 to 26 of 26

Filter...

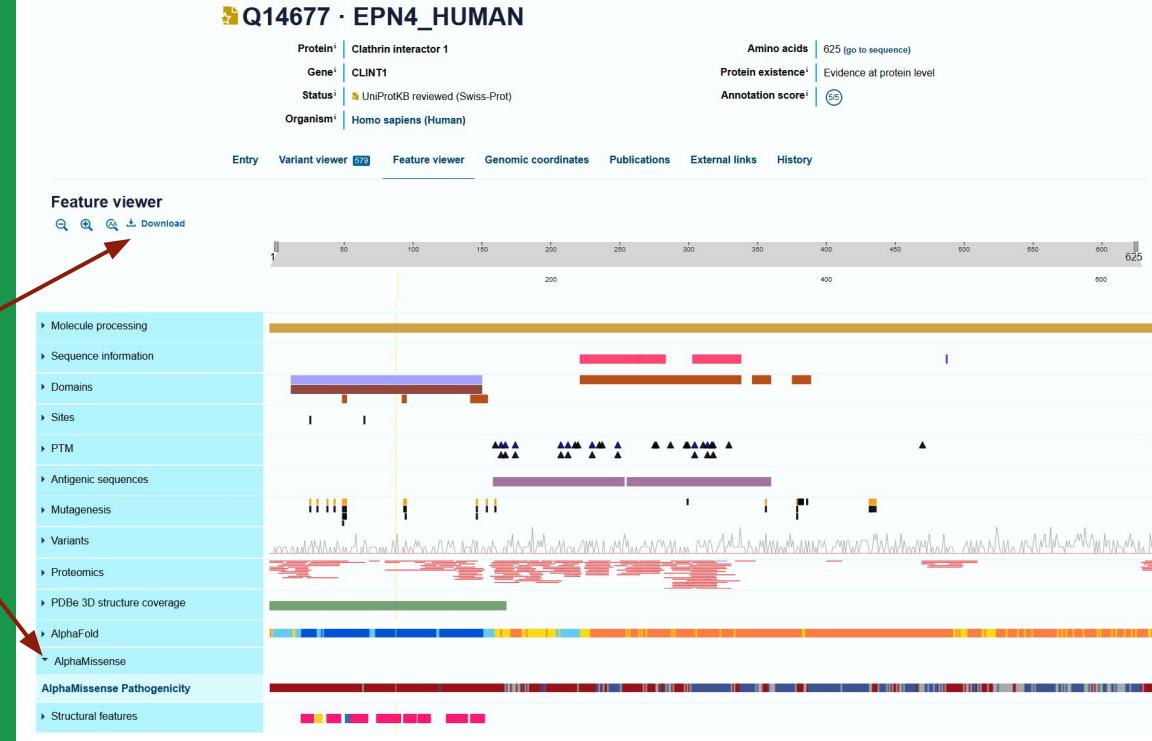
Transcript	Protein change	VEP Consequence	Other annotations
ENST00000544040.7 - NM_001273.5 <small>MANE Select Selected transcript</small> ENST00000544040.7:c.1933C>T	R/W at position 645 of 1912 ENSP0000440542.2:p.Arg645Trp	missense_variant	Sift PolyPhen CADD REVEL AlphaMissense SpliceAI phyloP Deleterious low confidence (0) Probably damaging (0.994) 32 0.838 Likely Pathogenic (0.9183) ≤ 0.2 2.765
ENST00000645022.1 NM_001297553.2 ENST00000645022.1:c.1912C>T	R/W at position 638 of 1905 ENSP0000496163.1:p.Arg638Trp	missense_variant	Sift PolyPhen CADD REVEL SpliceAI phyloP Deleterious low confidence (0) Probably damaging (0.967) 32 0.838 ≤ 0.2 2.765

AlphaMissense scores and classifications of likely pathogenicity are displayed on DECIPHER variant pages.

They are only for the transcript which matches the UniProt protein used in score calculation.

AlphaMissense in UniProt

- UniProt is a comprehensive, high-quality and freely accessible resource of protein sequence and functional information for a large variety of species.
- AM annotations can be downloaded by opening the Download panel.
- AlphaMissense scores are integrated in the feature viewer, allowing the visualization of overlaps of AM scores in context with protein features such as domains, motifs, PTMs or variants (*heatmap view coming soon*).



AlphaMissense in ProtVar

- ProtVar is a resource to investigate SNV missense variation in humans by presenting annotations which may be relevant to interpretation. Developed in the Protein Function Development group.
- AM predictions can be viewed for 1 or multiple variants on any given list of human proteins.
- AlphaMissense scores are shown alongside other pathogenicity prediction scores and other functional annotations from UniProt.

GENOMIC

Chr.	Coordinate	ID	Ref.	Alt.	Gene	Codon (strand)	CADD	Isoform	Protein name	AA pos.	AA change	Consequence(s)	AlphaMiss. pred.	Annotations		
5	157817494		G	A	CLINT1	aCg/aUg (-)	29.1	can Q14677	Clathrin interactor	32	Thr/Met	missense	1.00			
Processed 1 input (1 protein)																

PROTEIN

Functional information

Variant Residue Position

Annotations from UniProt
No functional data for the variant position

Region Containing Variant Position

Annotations from UniProt

- Chain-Clathrin interactor 1
- Functional Domain-ENTH
- Structure predictions^{ref}
- Pockets containing variant
- Protein-protein interfaces containing variant

Chemical structures showing the conversion from Threonine to Methionine:

Threonine → Methionine

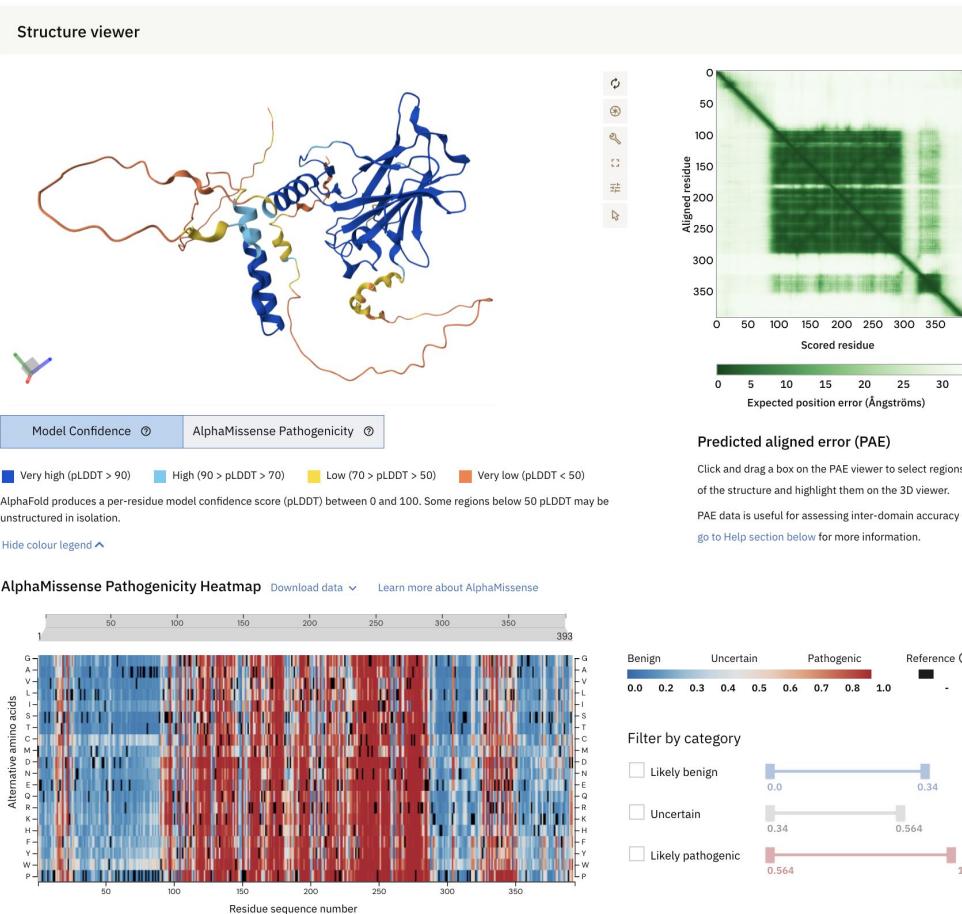
Pathogenicity predictions

Prediction Type	Score	Description
Conservation ^{ref}	0.99 i	very high i
Structure predictions		
Stability change $\Delta\Delta G$ ^{i ref}	1.97 i	unlikely to be destabilising ⁱ
Pathogenicity predictions		
CADD ^{ref}	29.1 i	probably deleterious ⁱ
AlphaMissense ^{ref}	1.00 i	pathogenic ⁱ
ESM-1b ^{ref}	-16.4 i	pathogenic ⁱ

ProtVar standardised coloursⁱ

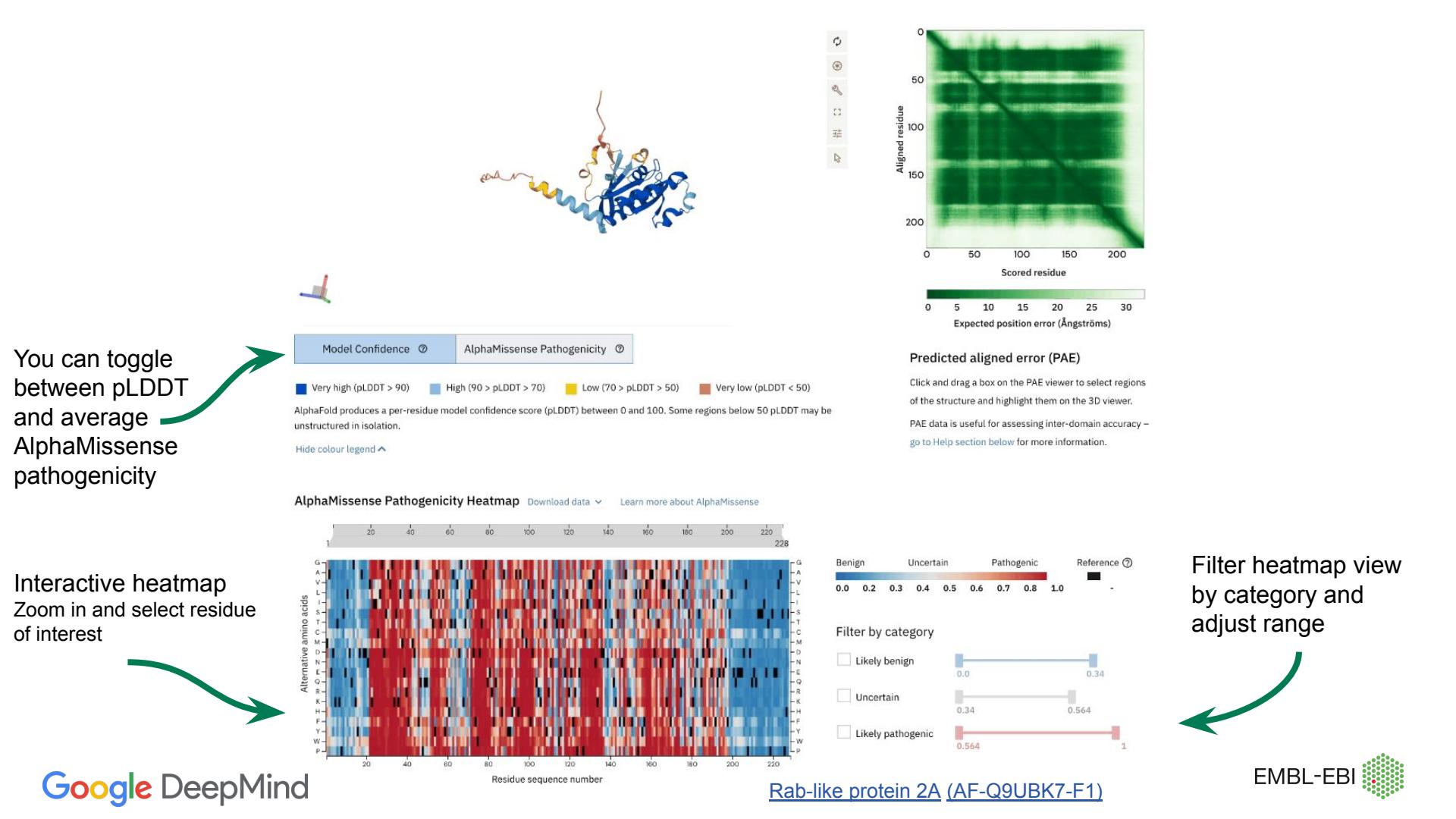
AlphaFold Protein Structure Database

Amino acid-level confidence measure, pLDDT, a local accuracy metric.



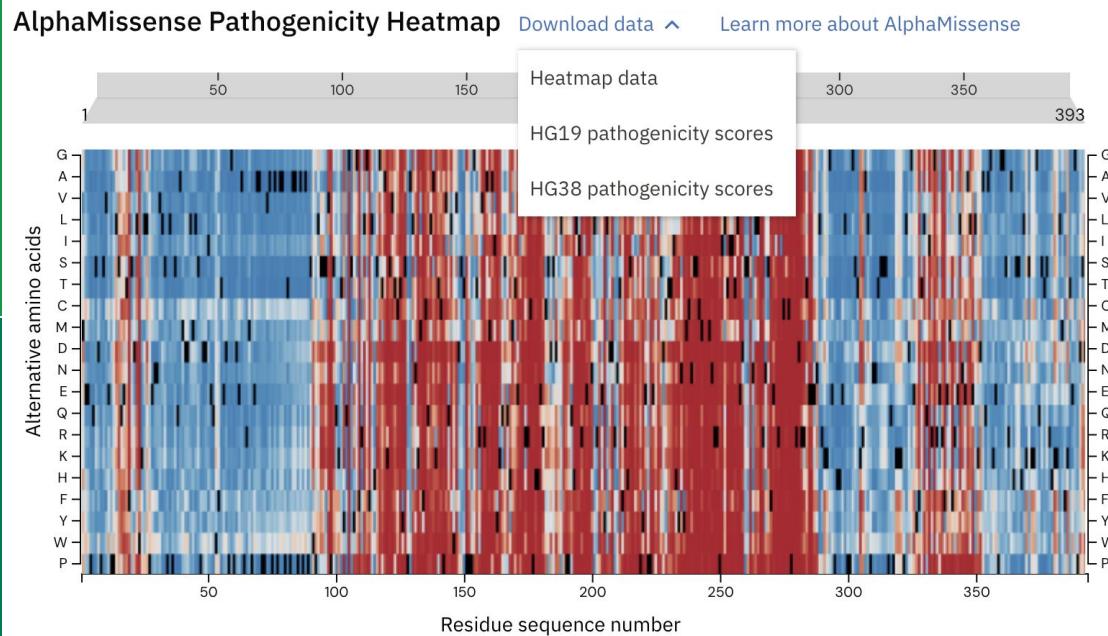
Measures the confidence in the relative position of two amino acids in Ångströms.

Integration of AlphaMissense data



Download AlphaMissense data

- **Heatmap data:** Contains the scores estimating the likelihood of pathogenicity and classifications for each possible amino acid substitution within the protein. **Used to visualise the heatmap you see on the entry pages.**
- **Pathogenicity scores** (HG19 and HG38). These files include AlphaMissense scores for all possible missense **single nucleotide variants across the human reference genome.** Each file corresponds to a specific genome assembly with information on specific genome position.



Download AlphaMissense data via API

```
▼ 0:
  entryId:          "AF-P04637-F1"
  gene:             "TP53"
  uniprotAccession: "P04637"
  uniprotId:        "P53_HUMAN"
  uniprotDescription: "Cellular tumor antigen p53"
  taxId:            9606
  organismScientificName: "Homo sapiens"
  uniprotStart:      1
  uniprotEnd:        393
  ▼ uniprotSequence: "MEEPQSDPSVEPLSQETFSDLWKLPPENNVLSPLSQAMDDMLSPDDIEQWFTEDPGPDEAPRMPEAAPPVAPAPAAPTP/"
  modelCreatedDate: "2022-06-01"
  latestVersion:    4
  ▼ allVersions:
    0:           1
    1:           2
    2:           3
    3:           4
    isReviewed:   true
    isReferenceProteome: true
  ▼ cifUrl:         "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-model\_v4.cif"
  ▼ bcifUrl:        "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-model\_v4.bcif"
  ▼ pdbUrl:         "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-model\_v4.pdb"
  ▼ paeImageUrl:   "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-predicted\_aligned\_error\_v4.png"
  ▼ paeDocUrl:     "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-predicted\_aligned\_error\_v4.json"
  ▼ amAnnotationsUrl: "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-aa-substitutions.csv"
  ▼ amAnnotationsHg19Url: "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-hg19.csv"
  ▼ amAnnotationsHg38Url: "https://alphafold.ebi.ac.uk/files/AF-P04637-F1-hg38.csv"
```



Response example for P04637:

<https://alphafold.ebi.ac.uk/api/prediction/P04637?key=AlzaSyCeurAJz7ZGjPQUtEaerUkBZ3TaBkXrY94>

Thank you!

Integration of AlphaMissense data into EMBL-EBI resources

7th June 2024

