# Project 1.3: Analyzing the Protein Coding Mutations in the Zimmerome
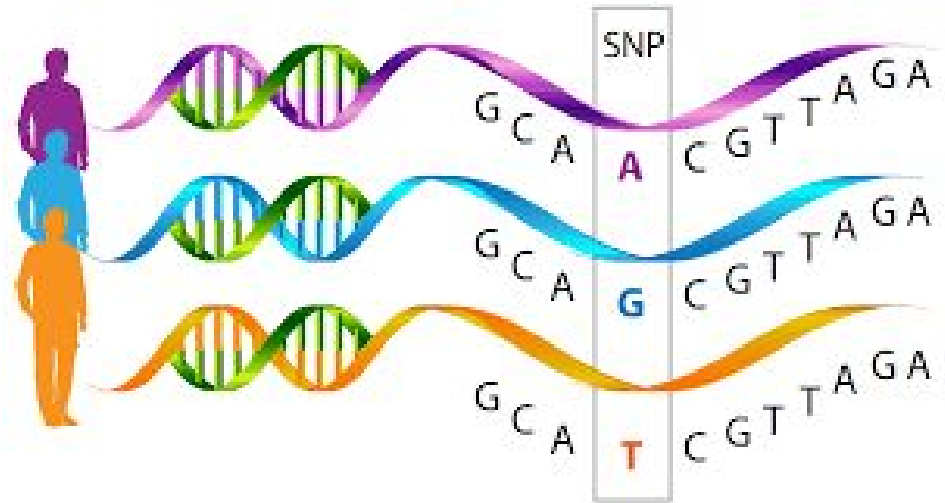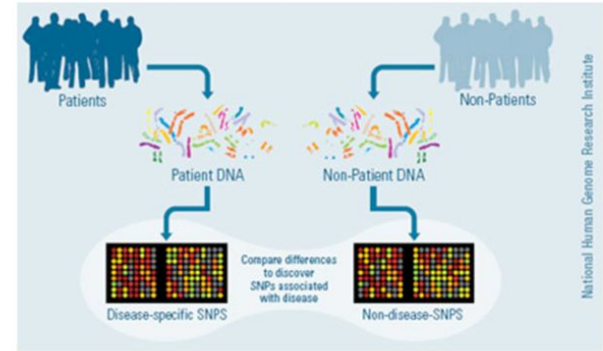
Ramya Prathuri, Megan Brady, and Nir Neumark

# What is Variant Prioritization?
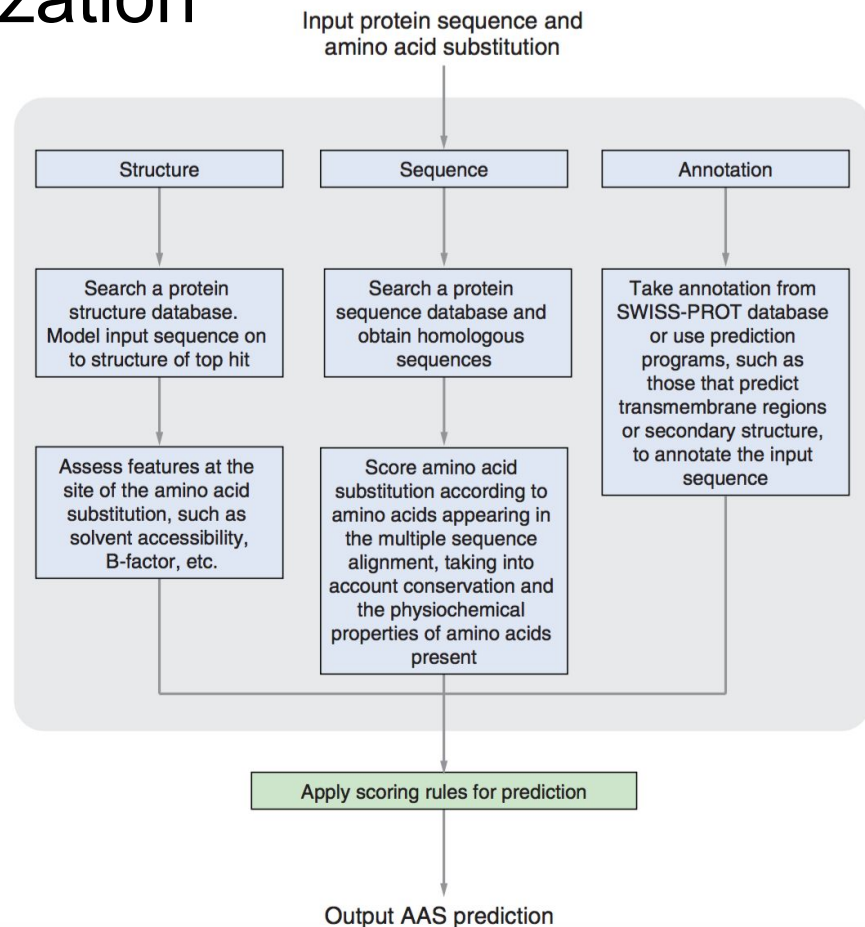
*The process of identifying deleterious SNVs*

GWAS



- SNP = seen in > 1% population
- SNV = no limitation on frequency

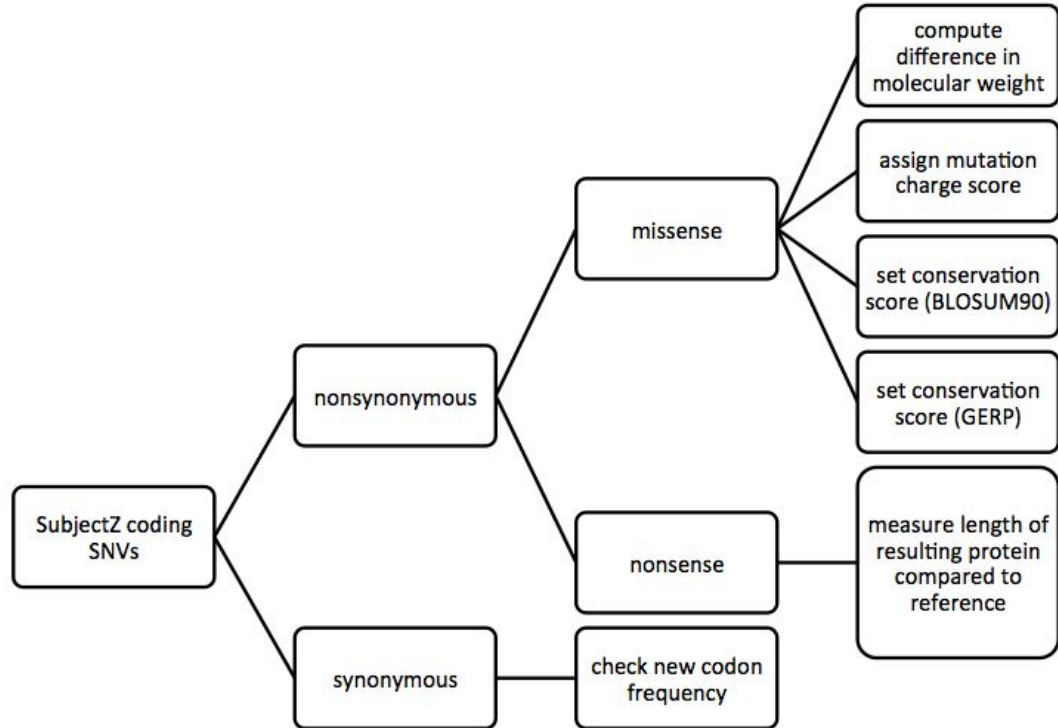# Principles of Variant Prioritization

- Sequence:
  - synonymous or not?
  - missense or nonsense?
  - aa charge? size?
- Structure:
  - Phyre2, TMHMM
- Annotation:
  - dsSNP? GO term?

Input protein sequence and amino acid substitution

| Structure | Sequence | Annotation |
|---|---|---|
| Search a protein structure database. Model input sequence on to structure of top hit | Search a protein sequence database and obtain homologous sequences | Take annotation from SWISS-PROT database or use prediction programs, such as those that predict transmembrane regions or secondary structure, to annotate the input sequence |
| Assess features at the site of the amino acid substitution, such as solvent accessibility, B-factor, etc. | Score amino acid substitution according to amino acids appearing in the multiple sequence alignment, taking into account conservation and the physiochemical properties of amino acids present | |

Apply scoring rules for prediction

Output AAS prediction

# Goal: How to prioritize variants?

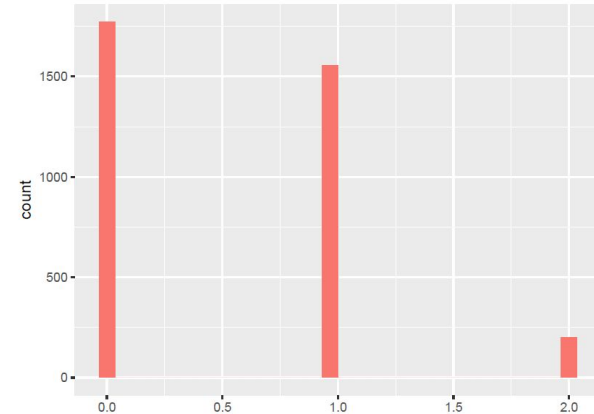*Consider a variety of parameters to assess SNVs*

# Mutation features

- Location
- Charge Change
- Size Difference
- Substitute
  - Blosum
  - PAM
  - GERD

# Charge Change

|            | Neutral | Hydrophilic | Hydrophobic |
|------------|---------|-------------|-------------|
| Neutral    | 0       | 1           | 1           |
| Hydrophilic| 1       | 0           | 2           |
| Hydrophobic| 1       | 2           | 0           |



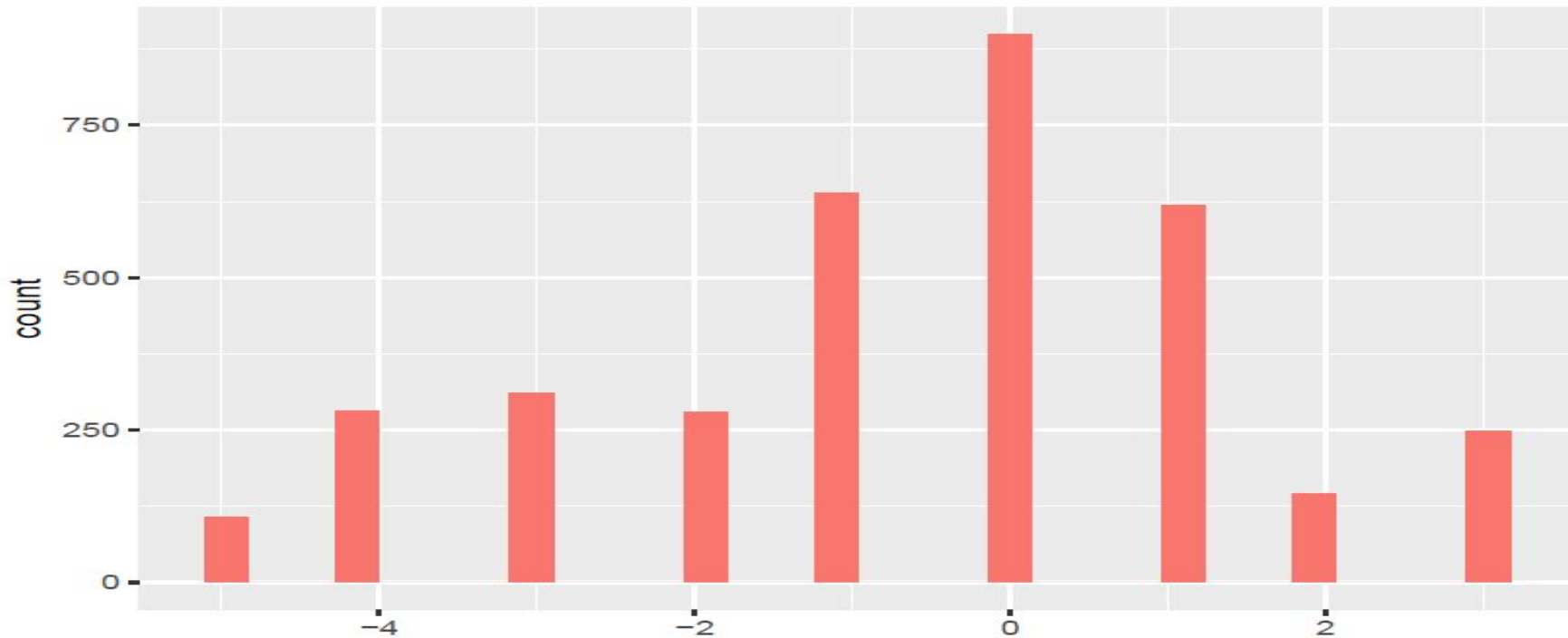| I | V | L | F | C | M | A | W [1] | G | T | S | Y | P | H | N | D | Q | E | K | R |
|---|---|---|---|---|---|---|-------|---|---|---|---|---|---|---|---|---|---|---|---|
| 4.5 | 4.2 | 3.8 | 2.8 | 2.5 | 1.9 | 1.8 | -0.9 | -0.4 | -0.7 | -0.8 | -1.3 | -1.6 | -3.2 | -3.5 | -3.5 | -3.5 | -3.5 | -3.9 | -4.5 |
| HYDROPHOBIC | | | | | | | | NEUTRAL | | | | | | HYDROPHILIC | | | | | |

# Size Difference

**Properties of amino acids**

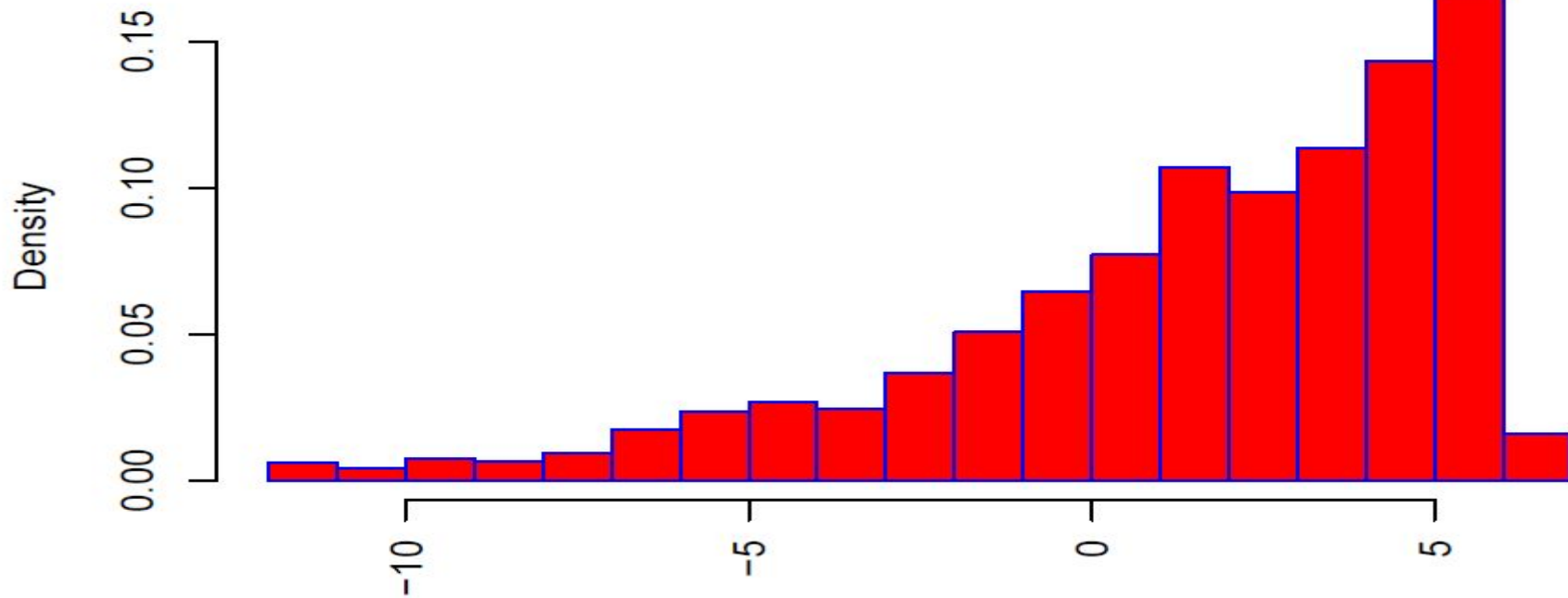| Amino acid residue | pKa of ionizing side chain[a] | Average residue mass[b] (daltons) | Monoisotopic mass (daltons)[b] | Occurrence in proteins[c] (%) | Percent buried residues[d] (%) | $V_r$[e] (Å³) | van der Waals volume[f] (Å³) | Accessible surface area[g] (Å²) | Ranking of amino acid polarities[h] |
|---|---|---|---|---|---|---|---|---|---|
| Alanine | – | 71.0788 | 71.03711 | 7.5 | 38 (12) | 92 | 67 | 67 | 9 (7) |
| Arginine | 12.5 (>12) | 156.1876 | 156.10111 | 5.2 | 0 | 225 | 148 | 196 | 15 (19) |
| Asparagine | – | 114.1039 | 114.04293 | 4.6 | 10 (2) | 135 | 96 | 113 | 16 (16) |
| Aspartic acid | 3.9 (4.4–4.6) | 115.0886 | 115.02694 | 5.2 | 14.5 (3) | 125 | 91 | 106 | 19 (18) |
| Cysteine | 8.3 (8.5–8.8) | 103.1448 | 103.00919 | 1.8 | 47 (3) | 106 | 86 | 104 | 7 (8) |
| Glutamine | – | 128.1308 | 128.05858 | 4.1 | 6.3 (2.2) | 161 | 114 | 144 | 17 (14) |
| Glutamic acid | 4.3 (4.4–4.6) | 129.1155 | 129.04259 | 6.3 | 20 (2) | 155 | 109 | 138 | 18 (17) |
| Glycine | – | 57.0520 | 57.02146 | 7.1 | 37 (10) | 66 | 48 | | 11 (9) |
| Histidine | 6.0 (6.5–7.0) | 137.1412 | 137.05891 | 2.2 | 19 (1.2) | 167 | 118 | 151 | 10 (13) |
| Isoleucine | – | 113.1595 | 113.08406 | 5.5 | 65 (12) | 169 | 124 | 140 | 1 (2) |
| Leucine | – | 113.1595 | 113.08406 | 9.1 | 41 (10) | 168 | 124 | 137 | 3 (1) |
| Lysine | 10.8 (10.0–10.2) | 128.1742 | 128.09496 | 5.8 | 4.2 (0.1) | 171 | 135 | 167 | 20 (15) |
| Methionine | – | 131.1986 | 131.04049 | 2.8 | 50 (2) | 171 | 124 | 160 | 5 (5) |
| Phenylalanine | – | 147.1766 | 147.06841 | 3.9 | 48 (5) | 203 | 135 | 175 | 2 (4) |
| Proline | – | 97.1167 | 97.05276 | 5.1 | 24 (3) | 129 | 90 | 105 | 13 (–) |
| Serine | – | 87.0782 | 87.03203 | 7.4 | 24 (8) | 99 | 73 | 80 | 14 (12) |
| Threonine | – | 101.1051 | 101.04768 | 6.0 | 25 (5.5) | 122 | 93 | 102 | 12 (11) |
| Tryptophan | – | 186.2133 | 186.07931 | 1.3 | 23 (1.5) | 240 | 163 | 217 | 6 (6) |
| Tyrosine | 10.9 (9.6–10.0) | 163.1760 | 163.06333 | 3.3 | 13 (2.2) | 203 | 141 | 187 | 8 (10) |
| Valine | – | 99.1326 | 99.06841 | 6.5 | 56 (15) | 142 | 105 | 117 | 4 (3) |

Average volume ($V_r$) of buried residues, calculated from the surface area of the side chain (Richards 1977; Baumann et al. 1989).

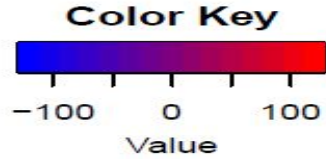# Blosum 90 Histogram

# GERP Score Histogram

# Mutation Frequency

| FI | LW | MR | CS | DV | EV | HL | HP | IF | LI | PQ | RI | SY | DY | GC | HN | LR | RM |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 |
| VD | Y* | CF | NY | QP | RL | SI | YF | LH | LQ | AG | EA | HD | MK | SF | YN | AD | FC |
| 3 | 3 | 4 | 4 | 4 | 4 | 4 | 4 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 7 | 7 |
| FY | PH | QK | SC | TK | DA | FV | NH | TR | IN | IS | LM | NK | RT | TN | VG | YH | AE |
| 7 | 7 | 7 | 7 | 7 | 8 | 8 | 8 | 8 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 10 |
| DH | WR | KT | ML | NT | QH | VF | YD | CW | IL | QL | EQ | GV | KQ | MI | PR | YC | AP |
| 10 | 10 | 11 | 11 | 11 | 11 | 11 | 11 | 12 | 12 | 14 | 15 | 15 | 15 | 15 | 15 | 15 | 16 |
| EG | LS | IM | AS | RP | CY | GD | KN | QE | HQ | DG | SL | SR | CG | CR | FS | RW | ST |
| 16 | 16 | 17 | 18 | 18 | 19 | 19 | 19 | 20 | 21 | 22 | 22 | 22 | 23 | 25 | 26 | 27 | 27 |
| EK | PA | GE | HY | LV | TP | GS | PT | RG | SA | GA | DN | FL | GR | LF | KE | RS | SP |
| 29 | 29 | 30 | 30 | 31 | 32 | 33 | 35 | 35 | 35 | 36 | 38 | 38 | 38 | 38 | 40 | 42 | 42 |
| DE | RK | ED | IT | KR | TS | PS | VL | LP | RC | SN | ND | NS | TI | MV | TM | VM | SG |
| 44 | 44 | 45 | 47 | 47 | 49 | 52 | 52 | 53 | 54 | 55 | 57 | 57 | 57 | 58 | 59 | 60 | 64 |
| HR | MT | AV | TA | QR | RQ | PL | RH | VA | VI | AT | IV | | | | | | |
| 65 | 66 | 70 | 75 | 81 | 85 | 88 | 89 | 97 | 108 | 111 | 129 | | | | | | |

# Mutation Frequency (matrix presentation)

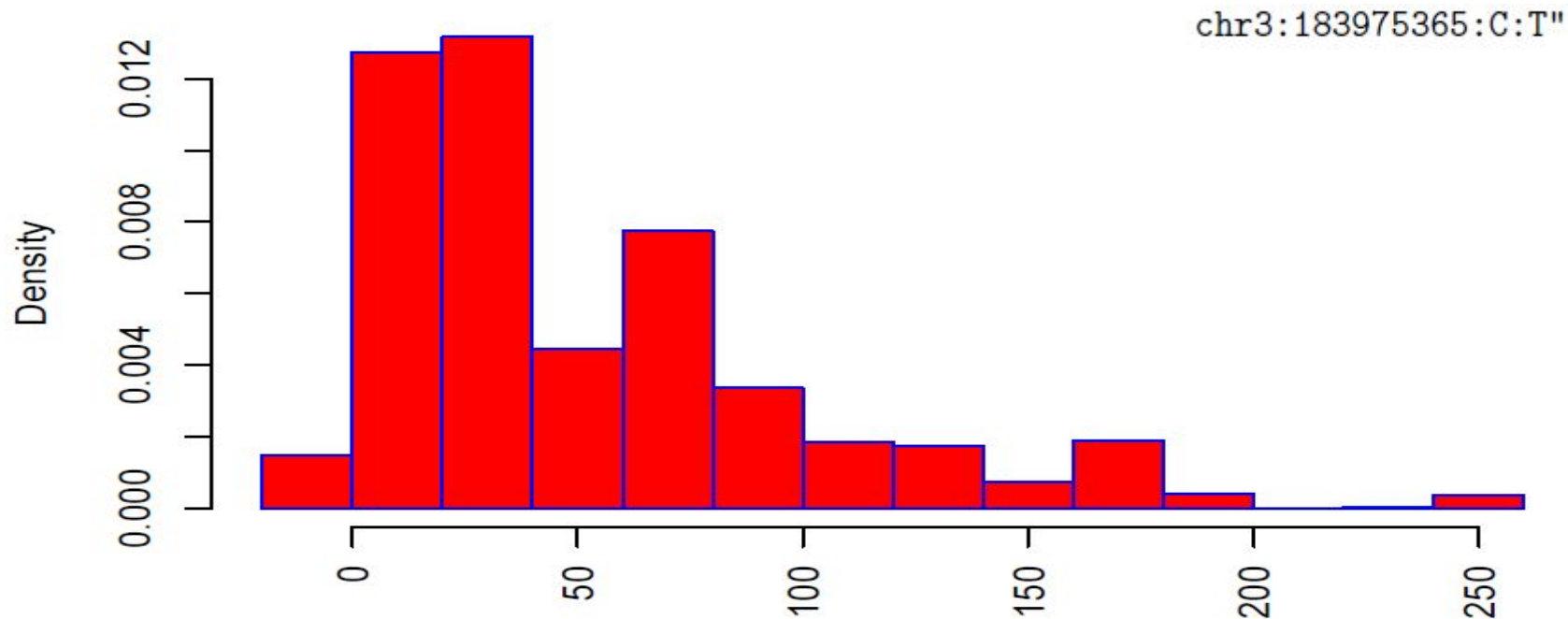|   | * | A | C | D | E | F | G | H | I | K | L | M | N | P | Q | R | S | T | V | W | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | O | O | O | 7 | 10 | O | 6 | O | O | O | O | O | O | 16 | O | O | 18 | 111 | 70 | O | O |
| C | O | O | O | O | O | 4 | 23 | O | O | O | O | O | O | O | O | 25 | 2 | O | O | 12 | 19 |
| D | O | 8 | O | O | 44 | O | 22 | 10 | O | O | O | O | 38 | O | O | O | O | O | 2 | O | 3 |
| E | O | 6 | O | 45 | O | O | 16 | O | O | 29 | O | O | O | O | 15 | O | O | O | 2 | O | O |
| F | O | O | 7 | O | O | O | O | O | 1 | O | 38 | O | O | O | O | O | 26 | O | 8 | O | 7 |
| G | O | 36 | 3 | 19 | 30 | O | O | O | O | O | O | O | O | O | O | 38 | 33 | O | 15 | O | O |
| H | O | O | O | 6 | O | O | O | O | O | O | 2 | O | 3 | 2 | 21 | 65 | O | O | O | O | 30 |
| I | O | O | O | O | O | 2 | O | O | O | O | 12 | 17 | 9 | O | O | O | 9 | 47 | 129 | O | O |
| K | O | O | O | O | 40 | O | O | O | O | O | O | O | 19 | O | 15 | 47 | O | 11 | O | O | O |
| L | O | O | O | O | O | 38 | O | 5 | 2 | O | O | 9 | O | 53 | 5 | 3 | 16 | O | 31 | 1 | O |
| M | O | O | O | O | O | O | O | O | 15 | 6 | 11 | O | O | O | O | 1 | O | 66 | 58 | O | O |
| N | O | O | O | 57 | O | O | O | 8 | O | 9 | O | O | O | O | O | O | 57 | 11 | O | O | 4 |
| P | O | 29 | O | O | O | O | O | 7 | O | O | 88 | O | O | O | 2 | 15 | 52 | 35 | O | O | O |
| Q | O | O | O | O | 20 | O | O | 11 | O | 7 | 14 | O | O | 4 | O | 81 | O | O | O | O | O |
| R | O | O | 54 | O | O | O | 35 | 89 | 2 | 44 | 4 | 3 | O | 18 | 85 | O | 42 | 9 | O | 27 | O |
| S | O | 35 | 7 | O | O | 6 | 64 | O | 4 | O | 22 | O | 55 | 42 | O | 22 | O | 27 | O | O | 2 |
| T | O | 75 | O | O | O | O | O | O | 57 | 7 | O | 59 | 9 | 32 | O | 8 | 49 | O | O | O | O |
| V | O | 97 | O | 3 | O | 11 | 9 | O | 108 | O | 52 | 60 | O | O | O | O | O | O | O | O | O |
| W | O | O | O | O | O | O | O | O | O | O | O | O | O | O | O | 10 | O | O | O | O | O |
| Y | 3 | O | 15 | 11 | O | 4 | O | 9 | O | O | O | O | 6 | O | O | O | O | O | O | O | O |

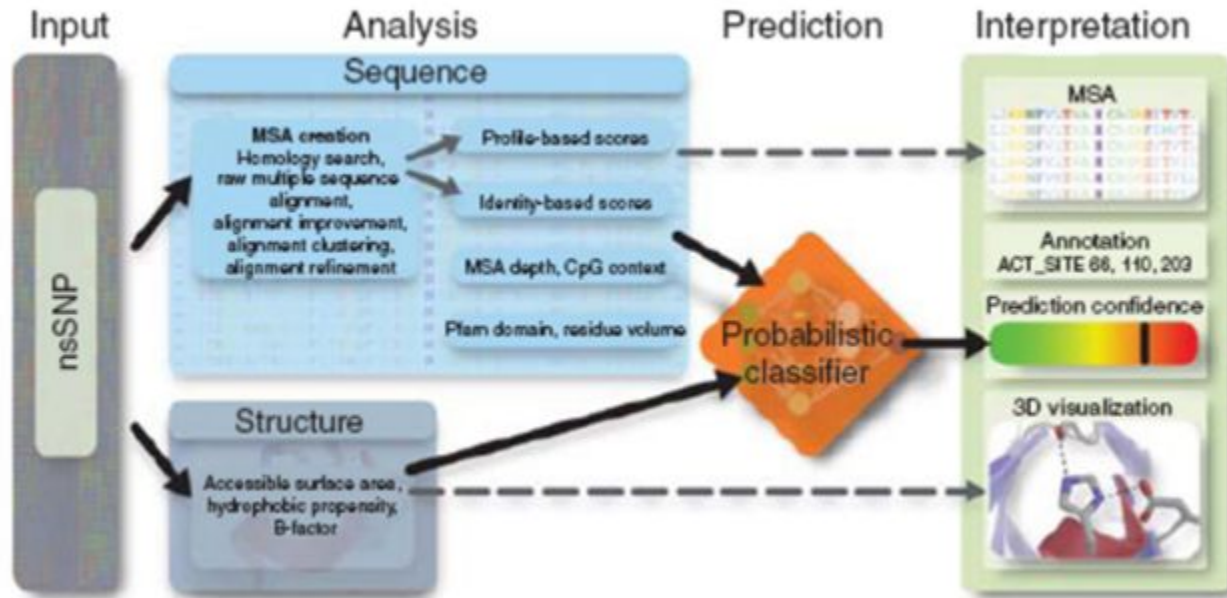# Mutation Frequency (heatmap)

# Deleterious Score (DScore)

$$DScore = \alpha Charge + \beta Size + \gamma GERP + \delta Blosum$$
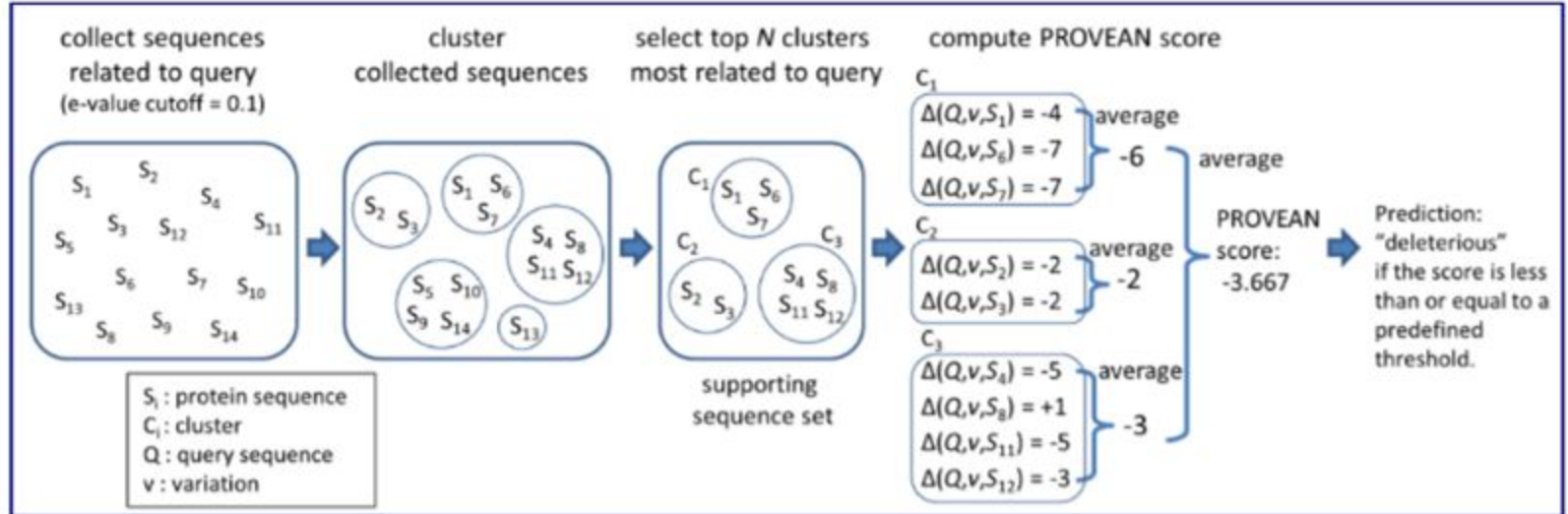
$$\alpha = \beta = \gamma = \delta = 1$$
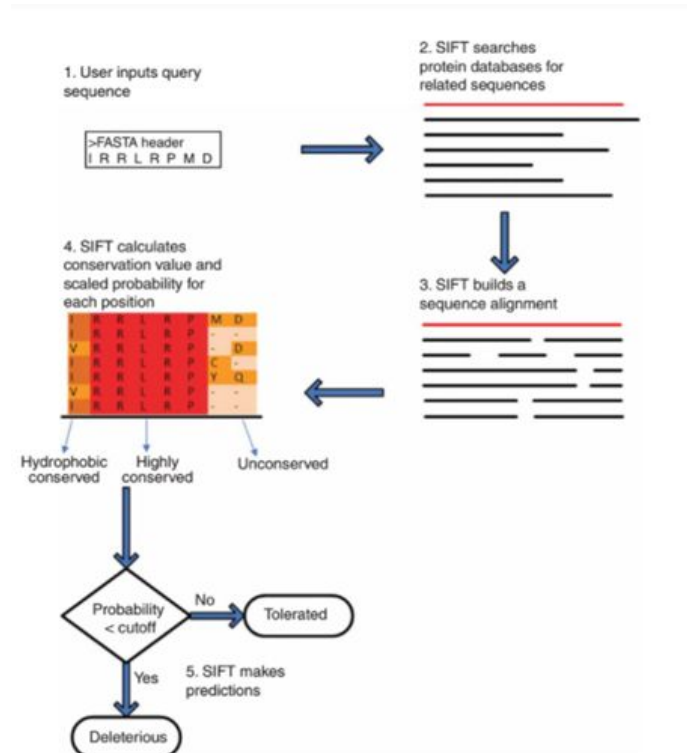
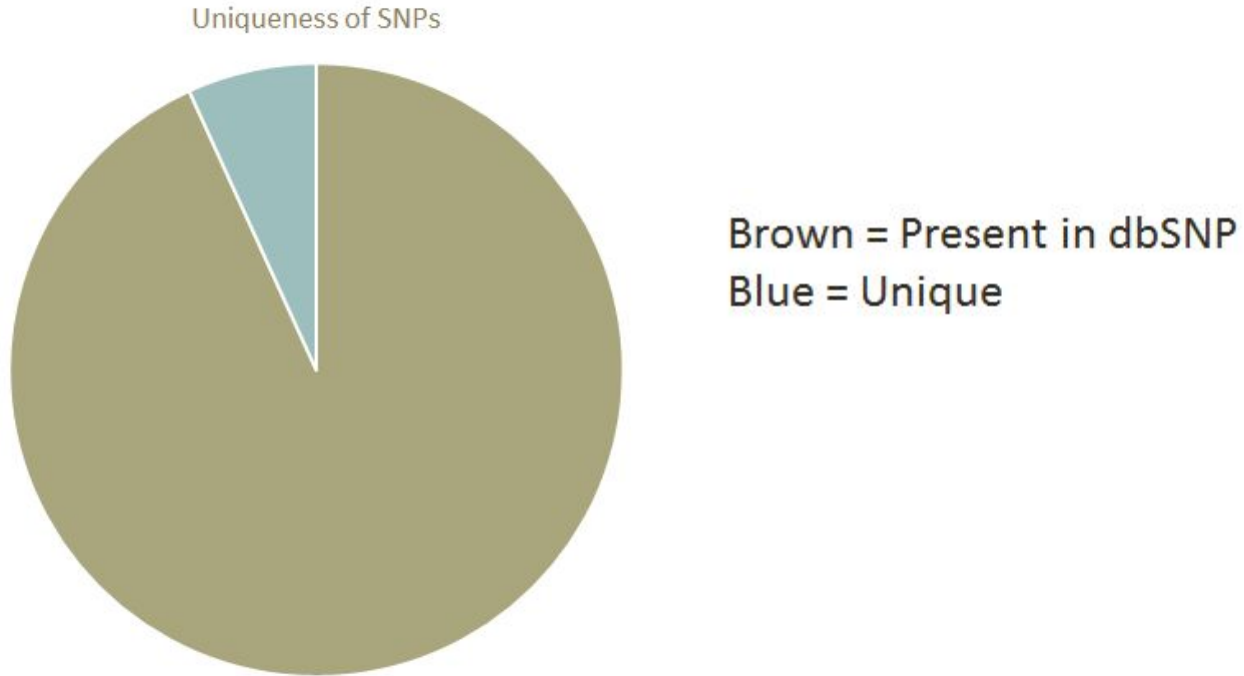# Deleterious Score Histogram

# PolyPhen

# PROVEAN

# SIFT

# PolyPhen2



- Probability calculated through Bayes classifier
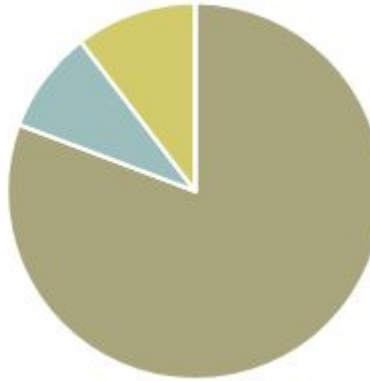- Benign (0-.45), possibly damaging (.45-.95), probably damaging (.95-1)

Adzhubei et al. 2015

# PROVEAN & SIFT

# Analyzing the Zimmer Genome

Uniqueness of SNPs



Brown = Present in dbSNP
Blue = Unique

# Comparison of Various Programs

|  | PolyPhen2 | PROVEAN | SIFT |
|---|---|---|---|
| Benign | 2669 | 3104 | 2851 |
| Possibly Damaging | 293 | 0 | 0 |
| Damaging | 344 | 420 | 658 |
| Undetermined | 0 | 6 | 21 |

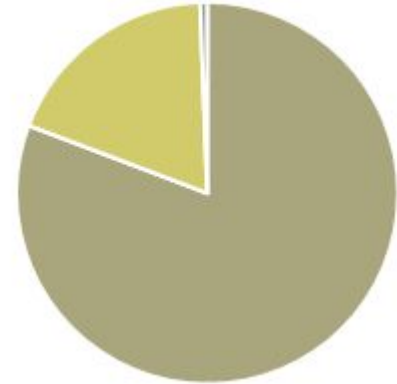|  | PolyPhen2 | PROVEAN | SIFT |
|---|---|---|---|
| Single Amino Acid Change | 3281 | 3514 | 3514 |
| Synonymous | 15 | 10 | 10 |
| Nonsense | 10 | 6 | 6 |

# Comparison of Various Programs



PolyPhen2 · PROVEAN · SIFT

Brown = Neutral
Yellow = Damaging
Blue = Possibly damaging
Gray = Undetermined

# Comparison of Various Programs



PROVEAN: 17,966 (90.3%)

543

853

659   1,146

15,618

153   457   469

SIFT: 16,887 (84.9%)   PolyPhen-2: 17,690 (88.9%)

19,898 disease variants

PROVEAN: 21,337 (61.5%)

1,454

6,784

2,528   1,111

16,244

2,184   2,991   1,405

SIFT: 23,947 (69.0%)   PolyPhen-2: 21,751 (62.7%)

34,701 common polymorphisms

http://provean.jcvi.org/about.php