

# Sequence Alignment-Based *In Silico* Pathogenicity Predictions for APC Missense Variants

Marc S. Greenblatt, M.D.

University of Vermont Larner College of Medicine



The University of Vermont  
LARNER COLLEGE OF MEDICINE



# Acknowledgements and Disclosures

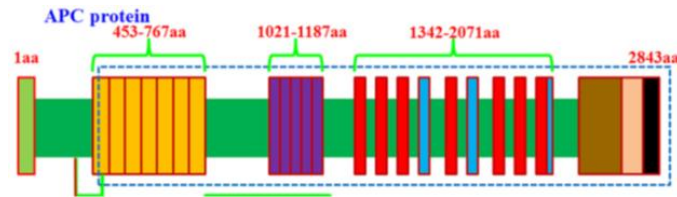
- University of Vermont
  - Alexander Karabachev
  - Dylan Martini
  - David Hermel
  - Dana Solcz
- Brown University
  - Indra Neil Sarkar
- Ambry Genetics
  - Tina Pesaran
  - Marcy Richardson



# Introduction

- *APC* variants causing FAP are loss of function, usually frameshift or nonsense
- Genetic testing frequently reveals *APC* missense substitutions, usually classified as Variants of Unknown Significance

Missense pathogenic variants are rare



Many functional domains are repeated

- *In silico* methods based on evolutionary sequence conservation are often used in hereditary cancer genes to help predict pathogenicity
- Unclear whether these *in silico* algorithms are useful in *APC* and other genes where non-missense mechanisms predominate- but they are often reported
- More data are needed to help interpret sequence-based methods for *APC* and less well-understood genes

# Computational Analysis

- *In silico* tools generate predictions of pathogenicity using:
  - Physiochemical characteristics of the substitution
  - Evolutionary conservation
  - Crystal structure if known
- Validated for BRCA1/2, MMR genes, TP53, not other genes
- Start with Protein Multiple Sequence Alignment (PMSA)

Sea Squirt	EYKVEDTPACFTPRSAISDLPCEEEDD
Sea Urchin	TYCVEGTPGPISRCSSLSSLDLNEELD
Zebrafish	TYCVEDTPICFSRGSSLSSLSSEEDM
Frog	TYCVEDTPICFSRGSSLSSLSAEDEI
Chicken	TYCVEDTPICFSRCSSLSSLSAEDEI
Opossum	TYCVEDTPICFSRCSSLSSLSAEDEI
Mouse	TYCVEDTPICFSRCSSLSSLSADDEI
Cow	TYCVEDTPICFSRCSSLSSLSAEDEV
Monkey	TYCVEDTPICFSRCSSLSSLSAEDEI
Human	TYCVEDTPICFSRCSSLSSLSAEDEI
	* **.* ** :: *::*.* ::

# PMSAs Require Much Curation

```
gi|699245605|ref|XP_009859498.1| TQTLKTVRSSLNGCLGVSSSLSDMLVAGSSDLLSVQHPVANDDTASMYSFST-----
gi|390338623|ref|XP_783363.3| KIPFGVAGSGNGGGSGAGSGNGE-----ETSSVMSFGSSSCSGTG
gi|219802769|ref|NP_001137312.1| ----GAAAAA-VCSQGSASRVDI-----DSASEMSSAGSY-----
gi|301613942|ref|XP_002936457.1| ----EIVTSSNVGSGQGSSSRADI-----DTASVMSSNSTY-----
gi|148236607|ref|NP_001084351.1| ----EITASGNVGSQGSSSRADI-----DTTSMVSSNSTY-----
gi|513231333|ref|XP_001233411.3| ----EISMSTS-NTGQGSAAARMDI-----ETASVMSSSNNY-----
gi|612035597|ref|XP_007497871.1| ----DVSVAPP-AGSQGSVAQVDO-----ETASGGGANGAY-----
gi|110225370|ref|NP_031488.2| ----ESNTAAS-SSGQSPATRVDI-----ETASVLSSSGTH-----
gi|274321915|ref|NP_001069454.2| ----EINMATS-GSGQGSTRIDI-----ETASVLSSSSTH-----
gi|297294842|ref|XP_002804524.1| ----EINMATS-GNGQGSTTRMDI-----ETASVLSSSSTH-----
gi|182397|gb|AAA03586.1| ----EINMATS-GNGQGSTTRMDI-----ETASVLSSSSTH----- 299
```

```
gi|699245605|ref|XP_009859498.1| ----SLPRQLAGNVSGSK-----TGEVCSL----LSSHDRHDMCTFQRLSQSEDSCI
gi|390338623|ref|XP_783363.3| TRKGSAAAGANNQQL-GVK-----VGFVYSLLSMLSSHDRDDMASTLLMMSRSADSC
gi|219802769|ref|NP_001137312.1| ----SVPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|301613942|ref|XP_002936457.1| ----SVPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|148236607|ref|NP_001084351.1| ----SVPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|513231333|ref|XP_001233411.3| ----SVPRLTSHL-GTKVTEYKQVEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|612035597|ref|XP_007497871.1| ----SVPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|110225370|ref|NP_031488.2| ----SAPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|274321915|ref|NP_001069454.2| ----SAPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|297294842|ref|XP_002804524.1| ----SAPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI
gi|182397|gb|AAA03586.1| ----SAPRLTSHL-GTK-----VEMVYSLLSMLGTHDKDDMSRTLLAMSSSQDSCI 346
* . :: **
```

Gaps or Insertions may represent:

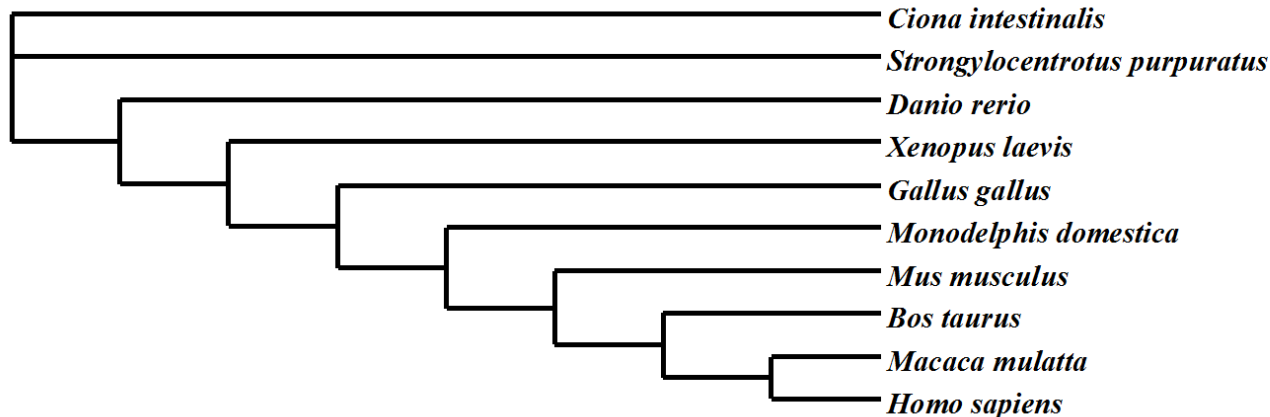
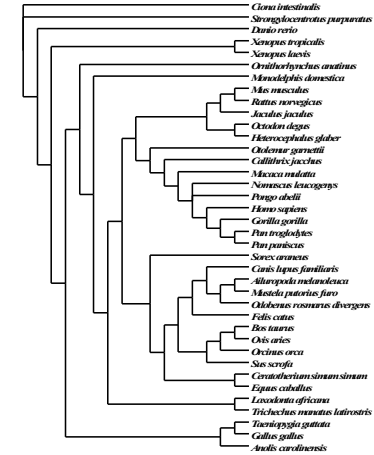
- True deletions and insertions
- Misreading of intron-exon boundaries
- Other sequence anomalies from genome assembly

# Our Goals

- Define the challenges of using *in silico* alignment-based tools for a protein with few pathologic missense variants.
- Create high quality PMSA for APC
- Generate predictions of pathogenicity for APC variants from *in silico* methods
- Compare to clinical classifications

# Methods for Phylogenetic Sequence Analysis To Validate Alignments Applied to APC

- Phylogenetic Trees constructed from 38 species, subsets
- $\geq 10$  species is enough APC evolutionary variation for statistically significant predictions (Greenblatt 2003, Cooper 2003)



# Curating the APC PMSA

- Created a PMSA using Clustal Omega
  - 10 species with evolutionary depth to sea squirt
  - Manually curated segments of the alignment (took some effort)
  - BLAST searches, removed exons that did not align with human
  - Located exon-intron boundaries, found omitted exons

Sea Squirt	SGRILTNLTYADNLN <b>KV</b> LLMNRGLLETVRDQLQHESEEIQ <b>DA</b> MASILRNLSWQADKEGR	479
Sea Urchin	AGMALTNLTFGDVTKALLCSMKGCMKALVALLSAESEDLRQVAASVLRNLSWRADMASK	502
Zebrafish	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRAMVAQLKSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	554
Frog	AGMALTNLTFGDVAN <b>KAT</b> LCSMKSCMRALVAQLKSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	562
Chicken	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRALVAQLKSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	560
Opossum	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRALVAQLKSESEDLE <b>QV</b> IASVLRNLSWRADVNSK	560
Mouse	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRALVAQLKSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	558
Cow	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRALVAQLQSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	560
Monkey	AGMALTNLTFGDVAN <b>KAT</b> LCSMKGCMRALVAQLKSESEDLO <b>QV</b> IASVLRNLSWRADVNSK	560
Human	<b>AGMALTNLTFGDVAN<b>KAT</b>LCSMKGCMRALVAQLKSESEDLO<b>QV</b>IASVLRNLSWRADVNSK</b>	560

.\* \*\*\*\*\*.\* \*\*.\* \*.\*. :.:. \*. \*\*\*\*\*.\* \*\*\*\*\*:\*\* ..

Sea Squirt	DLLRQKGVVRTLTECVICAKGEGTLKAMLSALWNLSGHCPANREDICSV <b>EG</b> SLAFLVSTL	539
Sea Urchin	KALREAGAVVALMTCSELEVKKES <b>TL</b> KS <sup>V</sup> LSALWNLSAHCTENKADICAVNGALEFLVSSL	562
Zebrafish	KILREVGSVRALMECALEV <b>QKE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICTVPGALAFVSTL	614
Frog	KTLREVGSVKALMECALDV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICVDGALAFVSTL	622
Chicken	KTLREVGSVKALMECALEV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICAVDGAFLVGT	620
Opossum	KTLREVGSVKALMECALEV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICAVDGAFLVGT	620
Mouse	KTLREVGSVKALMECALEV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICAVDGAFLVGT	618
Cow	KTLREVGSVKALMECALEV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICAVDGAFLVGT	620
Monkey	KTLREVGSVKALMECALEV <b>KE</b> STLKS <sup>V</sup> LSALWNLSAHCTENKADICAVDGAFLVGT	620
Human	<b>KTLREVGSVKALMECALEV<b>KE</b>STLKS<sup>V</sup>LSALWNLSAHCTENKADICAVDGAFLVGT</b>	620

. \*\*.\* \* \* : \* \* : . : \* .\*\*\*\*\*:\*\* \* : \*\*\*\*\* \* \* \* \* : \*



# APC Missense Variants Found in ClinVar

- From 2013 to 2018 the total number of reported APC missense variants increased from N= 47 → N= 1,988
- Missense variants in ClinVar, 5/2018, classifications:
  - Benign – 22 (1.1%)
  - Pathogenic – 9 (0.4%)
    - 2 somatic
    - 7 at splice sites → **not** “Likely Pathogenic” due to missense
  - Uncertain significance – 1806 **(93.2%)**
  - Conflicting interpretations of Pathogenicity – 103 **(5.3%)**
- **98.5% of APC missense classifications in ClinVar are not clinically useful.**

# *In silico* Predictions of Pathogenicity for APC Missense Variants

Method	ClinVar Benign Variants (N=22)	ClinVar Variants of Uncertain Significance (N=1903)
	Total Accuracy	Predicted Neutral (%)
<u>REVEL</u>	100%	N/A
<u>A-GVGD</u>	100%	82.5%
<u>SIFT</u>	95.5%	68.1%
<u>PolyPhen2</u>	81.8%	41.0%
<u>MAPP</u>	77.8%	25.0%

- Two likely pathogenic (LP) variants found in ClinVar p.S1028N, p.N1026S
- Classified as LP in July 2018 by Ambry Genetics using ACMG/AMP guidelines
- Evidence: protein features, segregation, published functional data
- Both located in the first 15-amino acid repeat of the  $\beta$ -catenin binding domain

# Conclusions

- Creating a high quality PMSA is labor intensive and a limiting factor in using *in silico* predictive tools
- *In silico* methods that are excellent classifiers for variants of some hereditary cancer genes (BRCA1/2, MMR, TP53) are not as accurate when applied to APC and some other genes
- Some APC features may predict why *in silico* methods perform poorly; these features may be common in other genes
- Systematic study of these features can improve predictive algorithms, interpretation of inherited genetic variation

