

Project Motivation

- Motivation: gain a greater understanding of mutations in ASD that will lead to better diagnostics
- What we know now: mutations in ASD probands are more likely to land in protein interaction domains than in their unaffected siblings
- Knowledge gap: the impact of multiple mutations on disrupting protein interactions has
 not been studied
- Model deployed and advantages: deep learning models allow for the effect of multiple mutations to be modeled in protein interactions
- Hypothesis: multiple mutations affecting protein interaction will be more common in ASD probands versus their unaffected siblings

Missense Variant Pipeline

OLD:

- Step 1: Pass two gene id strings to top layer function
- Step 2: Solicit/clean gene variant information
- Step 3: Map related proteins
- Step 4: Map protein sequences
- Step 5: Apply variant changes to sequences
- Step 6: Determine proband and sibling variants
- Step 7: Return new data frame for modeling*

NEW:

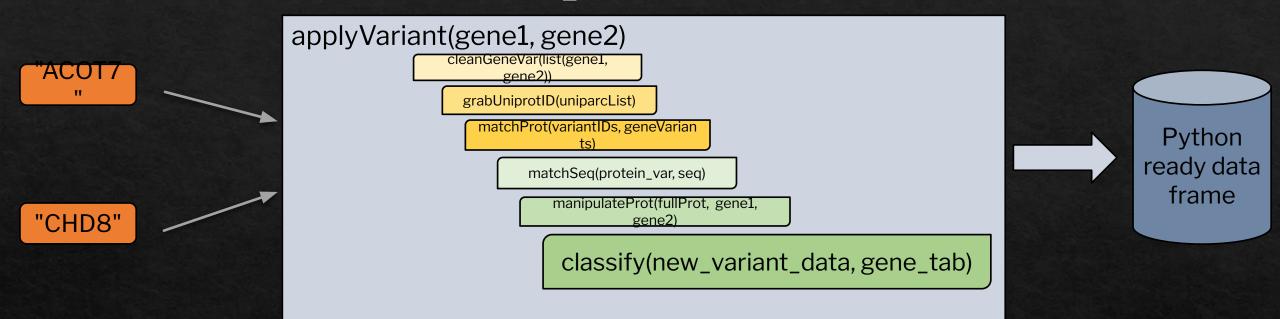
Step 1: Pass two gene id strings to top layer function

- Step 1a: Solicit/clean gene variant information
- Step 1b: Map related proteins
- Step 1c: Map protein sequences
- Step 1d: Apply variant changes to sequences (Keep Sequence string)
- Step 1e: Determine proband and sibling variants by number of occurrences
- Step 1f: Return new data frame for modeling/encoding

Step 2: Read file into python function

 Step 2a: Python function encodes and models data, returns Wilcoxon statistics

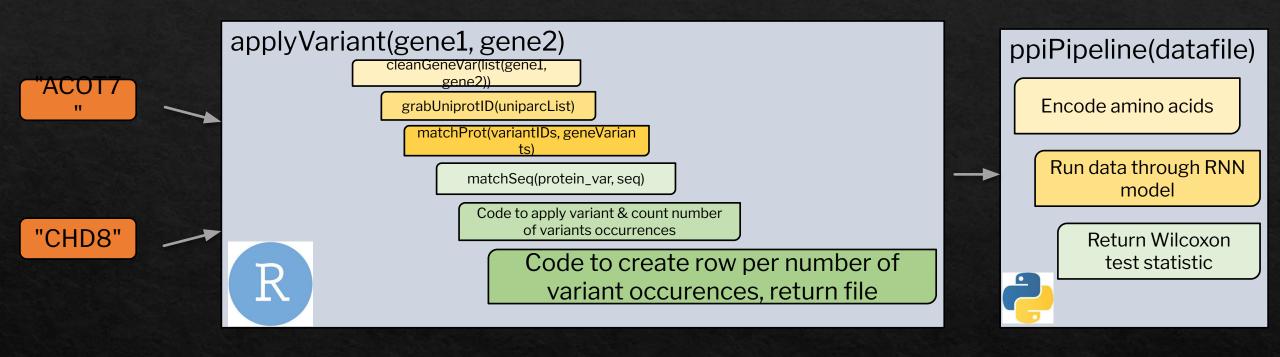
Last Time: Pipeline Architecture





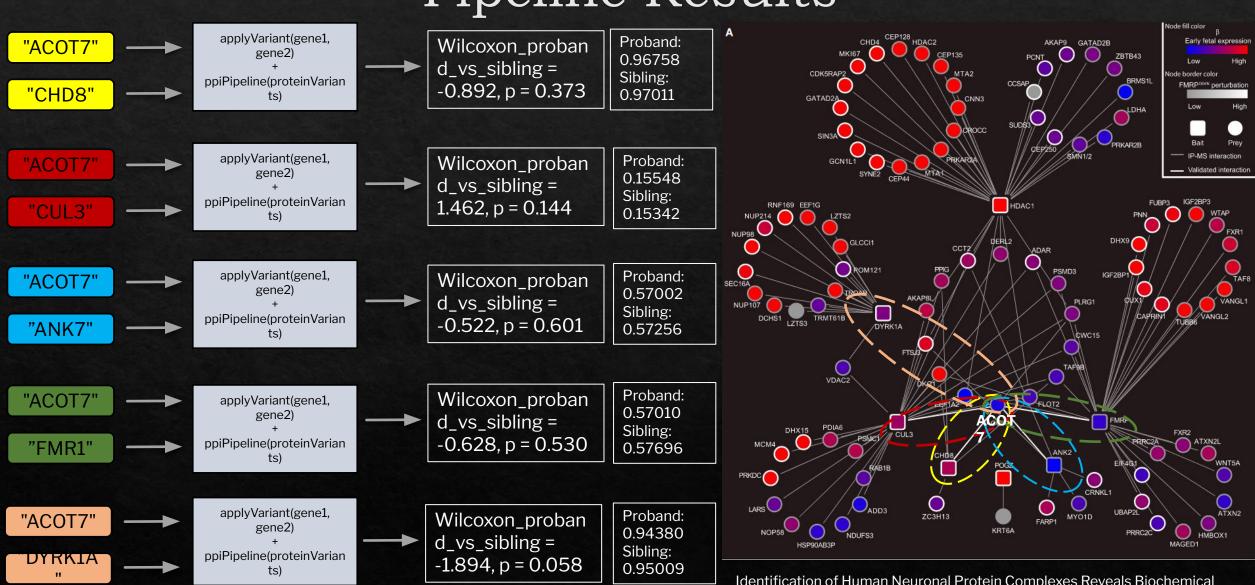
		tout						Carried Days of the Column of										THE RESIDENCE OF THE PARTY OF T	AND THE REAL PROPERTY OF THE PARTY OF THE PA	
	A	В	C	D	E	F	ú	Н	1000)	N I	La La	M	N	0	Ρ	Q	R	2	
1	Protein	Sequence	uniparc	symbol	gene	feature	biotype	havsp	protein_pos	seq_length	Old_aa	New_aa	Old_codon	New_codon	ensp	used_ref	variant_id	seq_matrix	variantSeq_matrix	source
2	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Asn227Ser	227	246	N	S	aAc	aGc	ENSP00000402532	T	36496	5 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	both
3	O00154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Lys225Asn	225	246	K	N	aaG	aaC	ENSP00000402532	C	36497	7 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	both
4	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Gly217Arg	217	246	G	R	Ggg	Agg	ENSP00000402532	C	36498	3 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	proband
5	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Met212Leu	212	246	M	L	Atg	Ctg	ENSP00000402532	T	36501	1 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	both
6	000154	MKLLARALRLCE	G• UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Gly205Arg	205	246	G	R	Gga	Aga	ENSP00000402532	C	36610	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	proband
7	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Pro172Leu	172	246	Р	L	cCa	<u>cTa</u>	ENSP00000402532	G	36627	7 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	proband
8	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Asp168His	168	246	D	Н	Gac	Cac	ENSP00000402532	C	36628	3 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	both
9	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Gly167Glu	167	246	G	E	gGg	gAg	ENSP00000402532	C	36629	9 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	proband
10	O00154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Tyr152Ter	152	246	Υ	*	taT	taG	ENSP00000402532	Α	36630	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	sib
11	O00154	MKLLARALRLCE	G-UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Arg149Gln	149	246	R	Q	cGg	cAg	ENSP00000402532	С	36633	3 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	sib
12	000154	MKLLARALRLCE	G UPI000002A4AF	ACOT7	ENSG00000097021	ENST00000418124	nonsense_mediate	ENSP00000402532.1:p.Glu146Asp	146	246	E) 9	D	gaG	gaC	ENSP00000402532	С	36634	4 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,	0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0	proband

Updated: Pipeline Architecture



4	A	В	С	D	E	F	G	Н	1	J	K	L	M	N	0	P	Q	R	S	T	U	
1	Protein	Sequence	uniparc	symbol	gene	feature	biotype	hgvsp	protein_pos	seq_length	Old_aa	New_aa	Old_codon	New_codon	ensp	used_ref	variant_id	seq_matrix	variantSeq_n	Proband	Sibling	
2	000154	MKLLARALR	UPI000002#	ACOT7	ENSG00000	ENST000006	protein_cod	ENSP000004	337	338	Q	K	Cag	Aag	ENSP000004	G	36411	MKLLARALR	MKLLARALR	1	,	0
3	000154	MKLLARALR	UPI000002#	ACOT7	ENSG00000	ENST000006	protein_cod	ENSP000004	334	338	Α	V	gCg	gTg	ENSP000004	G	36413	MKLLARALR	MKLLARALR	1	L	0
4	000154	MKLLARALR	UPI000002#	ACOT7	ENSG00000	ENST000006	protein_cod	ENSP000004	334	338	Α	V	gCg	gTg	ENSP000004	G	36413	MKLLARALR	MKLLARALR	1	L	0
5	000154	MKLLARALR	UPI000002#	ACOT7	ENSG00000	ENST000006	protein_cod	ENSP000004	331	338	Q	R	cAg	cGg	ENSP000004	T	36415	MKLLARALR	MKLLARALR	1		0
6	000154	MKLLARALR	UPI000002#	ACOT7	ENSG00000	ENST000006	protein_cod	ENSP000004	331	338	Q	R	cAg	cGg	ENSP000004	T	36415	MKLLARALR	MKLLARALR	1		0

Pipeline Results



Identification of Human Neuronal Protein Complexes Reveals Biochemical Activities and Convergent Mechanisms of Action in Autism Spectrum Disorders. Li, et al. *Cell Systems*, 2015.