

## Supporting information

### Multiscale simulations reveal architecture of NOTCH protein and ligand specific features

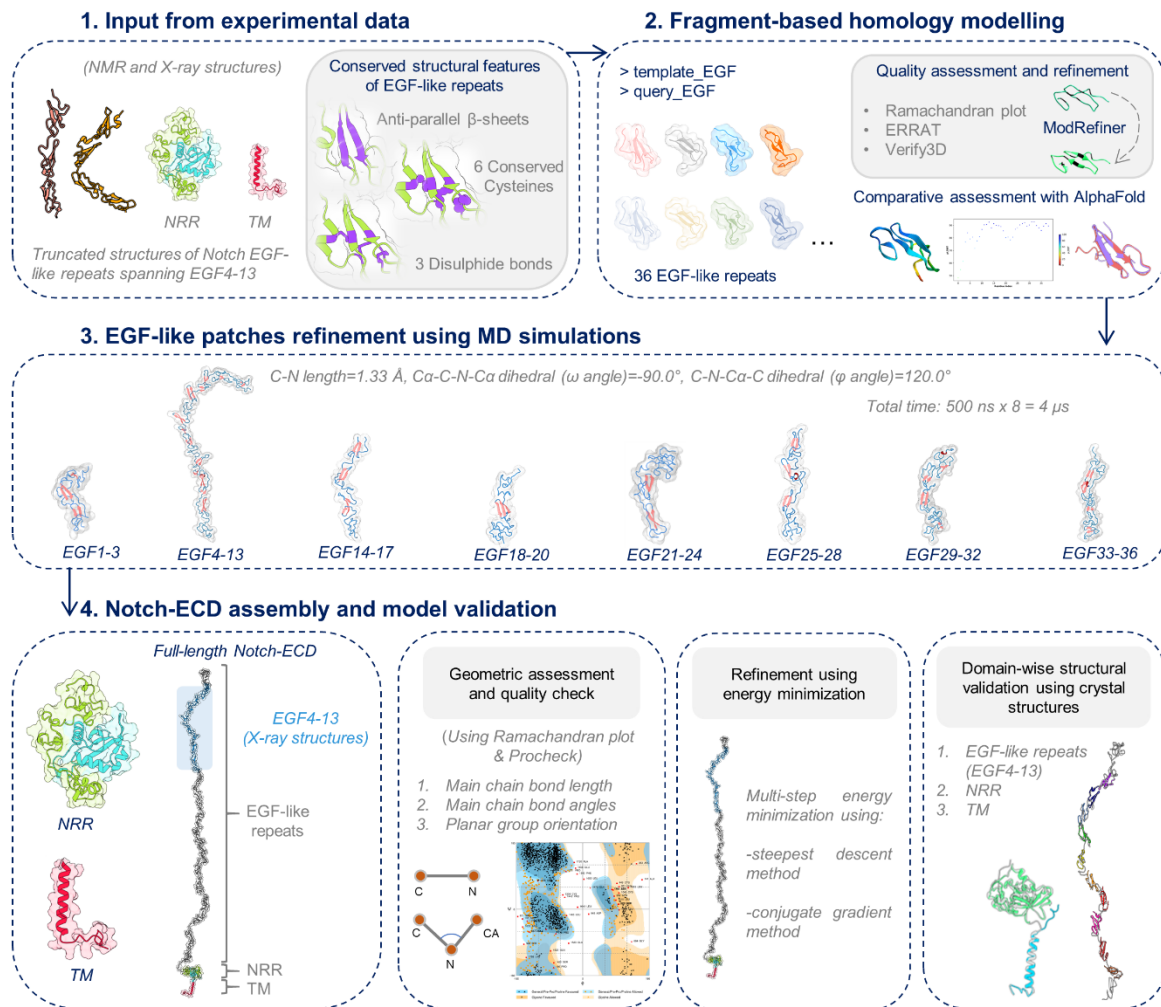
Surabhi Rathore<sup>1, 2</sup>, Deepanshi Gahlot<sup>1, 2</sup>, Jesu Castin<sup>1</sup>, Arastu Pandey<sup>3</sup>, Shreyas Arvindekar<sup>3</sup>, Shruthi Viswanath<sup>3</sup>, Lipi Thukral<sup>1, 2, \*</sup>

<sup>1</sup>CSIR-Institute of Genomics and Integrative Biology, Mathura Road, New Delhi- 110 025, India.

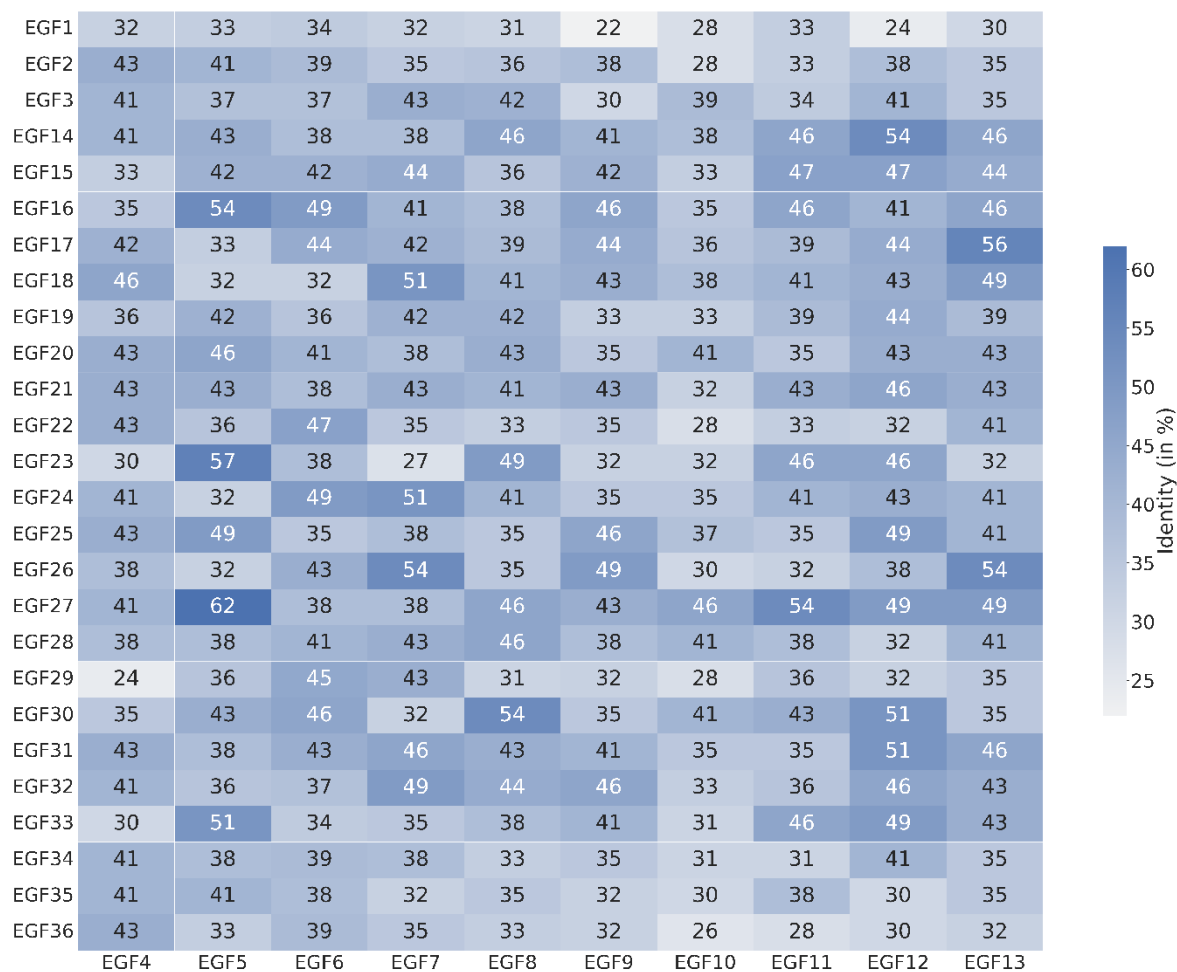
<sup>2</sup>Academy of Scientific and Innovative Research (AcSIR), Ghaziabad- 201002, India

<sup>3</sup>National Center for Biological Sciences (NCBS), Tata Institute of Fundamental Research (TIFR), GKVK, Bellary Road, Bangalore-560065, India

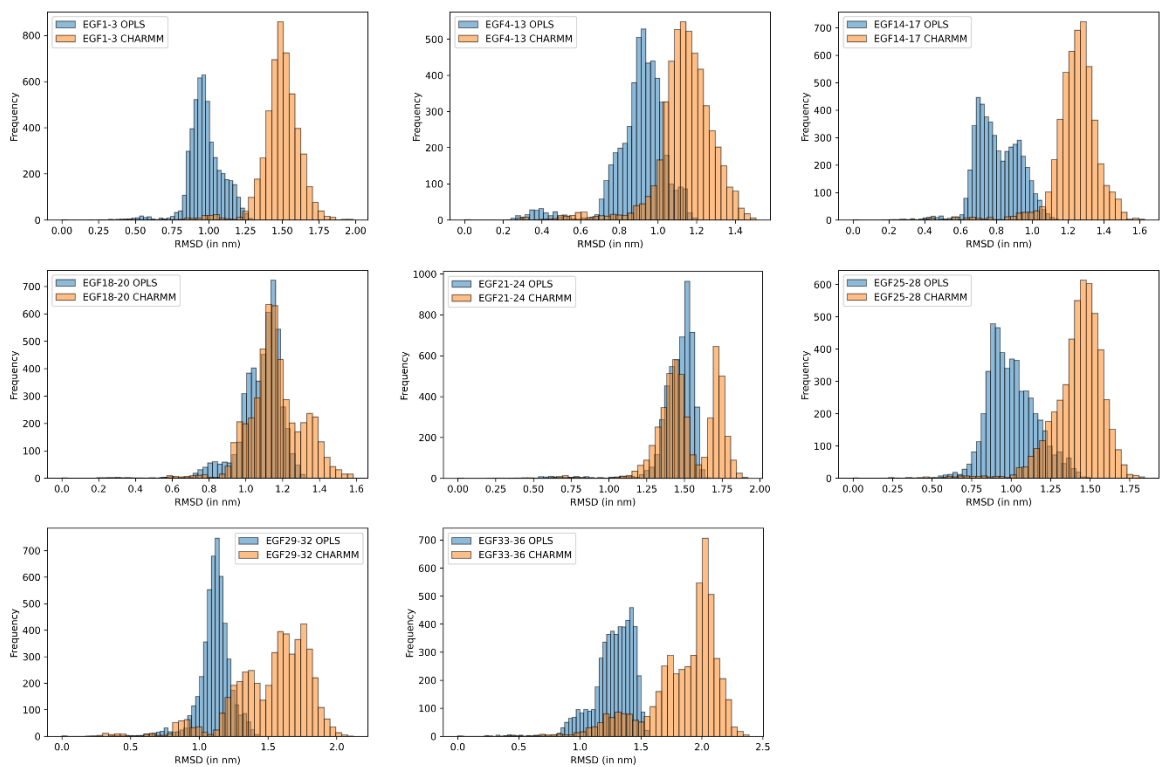
\*Correspondence to Lipi Thukral ([lipi.thukral@igib.res.in](mailto:lipi.thukral@igib.res.in))



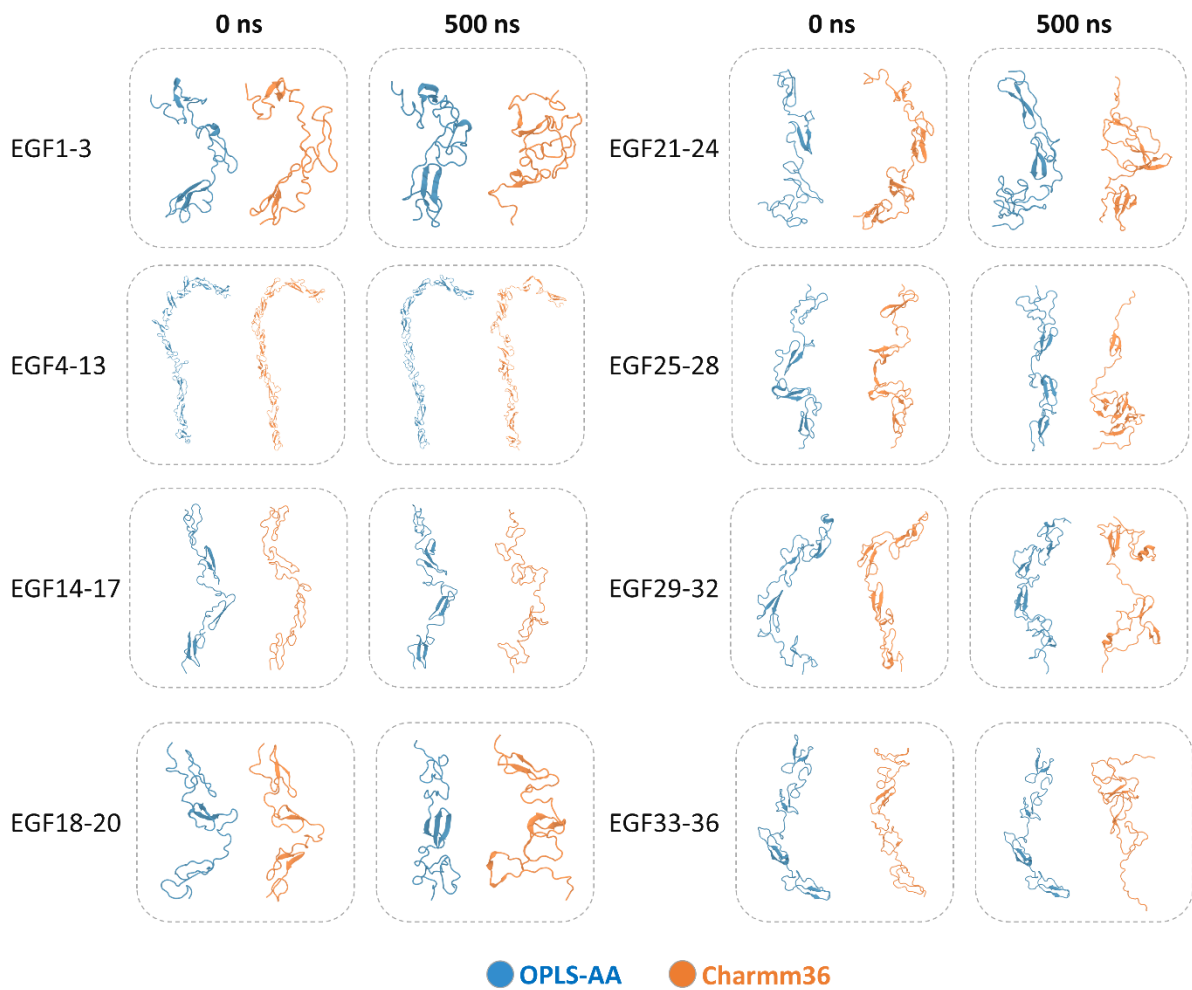
**Figure S1: Work approach pipeline to model NOTCH-ECD using integrative modelling approaches.** Step-wise detailed procedure employed for NOTCH modelling is given where input data was collected from NMR and X-ray crystal structures and we utilised the conserved features of EGF like repeats. The templates were selected based on the homology score between 36 EGFs of NOTCH ECD. The quality and refinement assessment were done followed by the MD simulations of the eight patches for further refinement. The eight patches were joined fulfilling the peptide bond criteria. The model was then subjected to various quality checks energy minimization refinement.



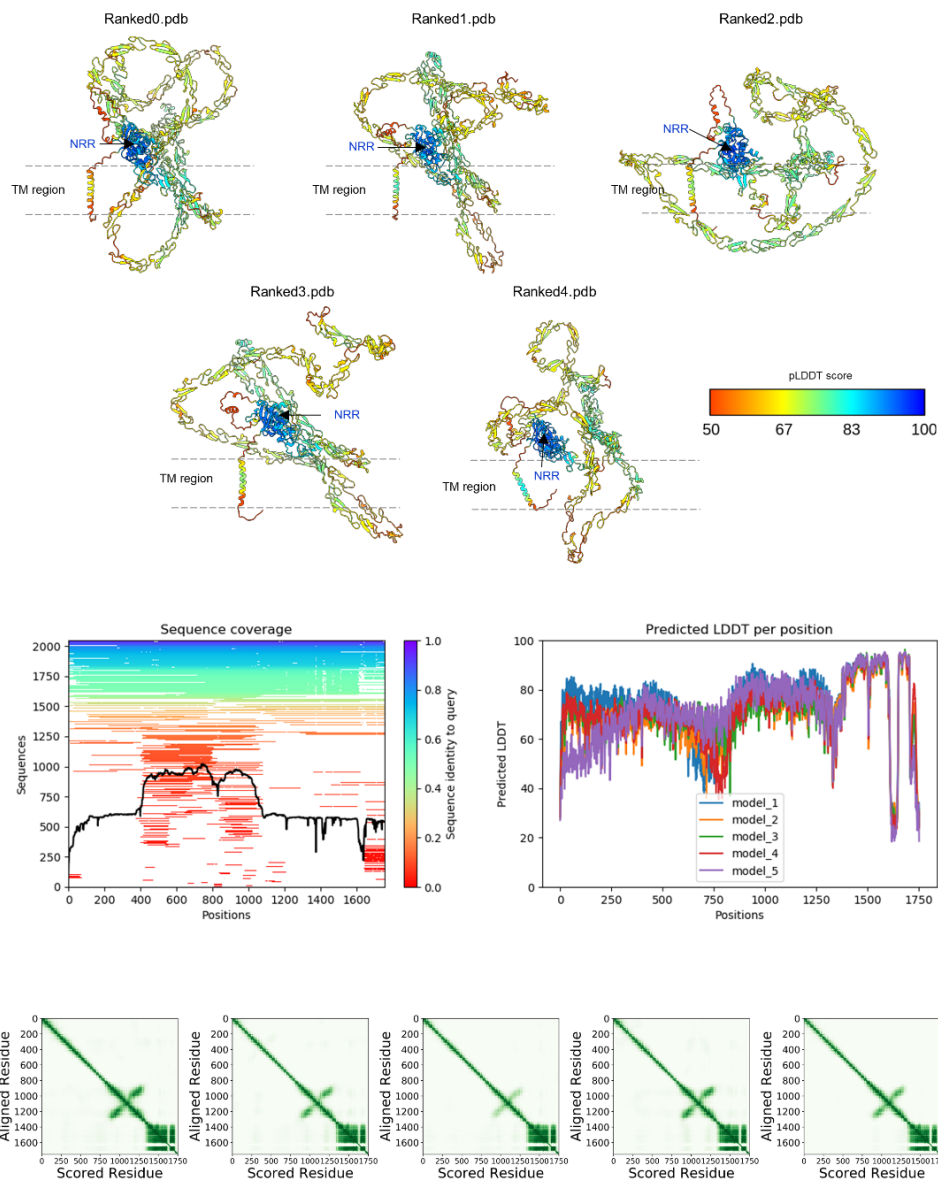
**Figure S2: Sequence identity score matrix for EGF-like repeats:** The sequence identity scores (in %) for EGF1-3 and EGF14-36 as compared to the EGF-like repeats with known crystal structures, i.e., EGF4-13 is shown in the heatmap, where the white to dark blue color gradient depicts the increasing identity.



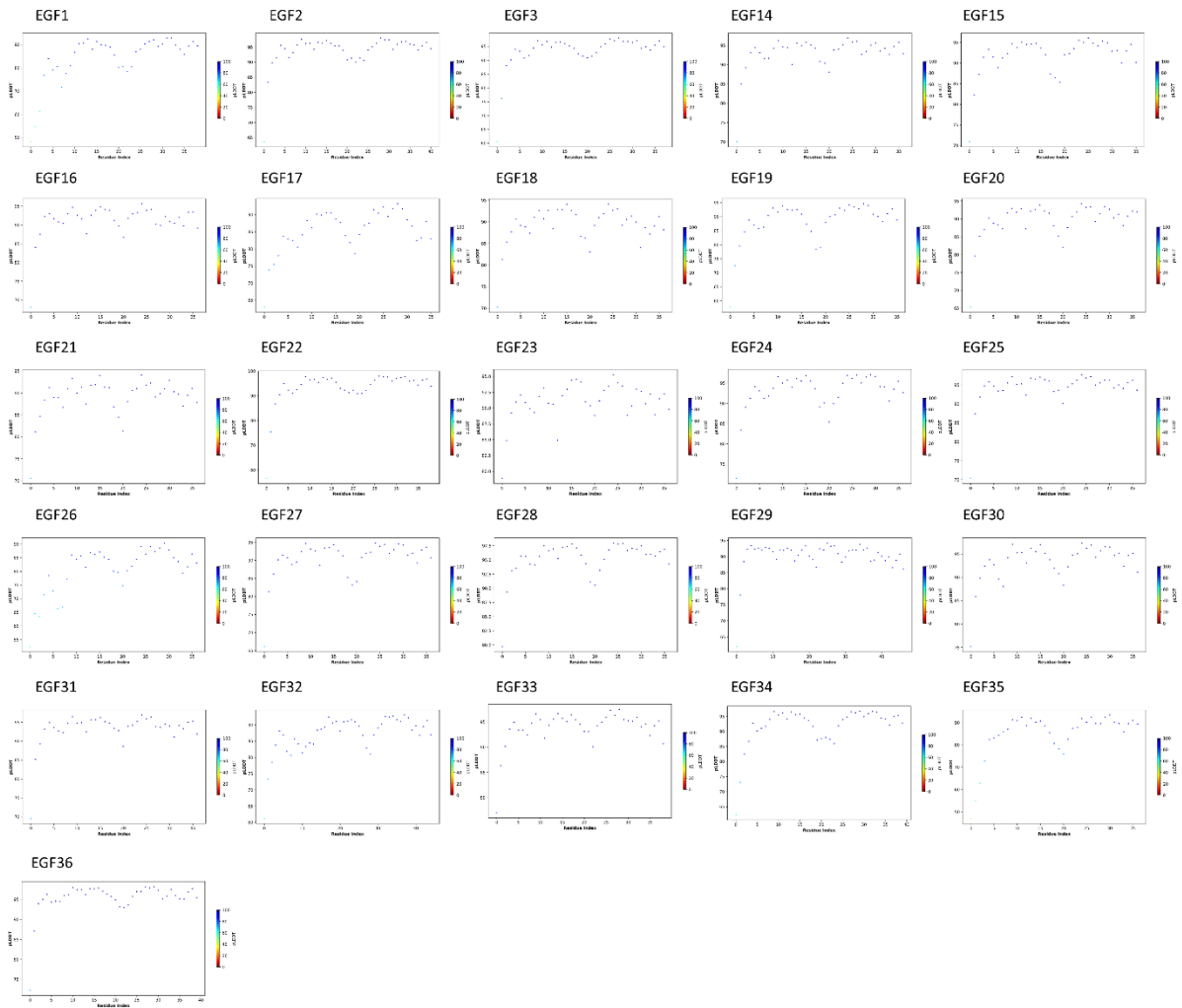
**Figure S3: Comparative assessment of EGF-like repeat patches refinement using OPLS and Charmm36:** The histogram displays a comparative root mean square deviation for each EGF-like repeat patch (EGF1-3, 4-13, 14-17, 18-20, 21-24, 25-28, 29-32, 33-36) refined using OPLS and Charmm36 force-fields.



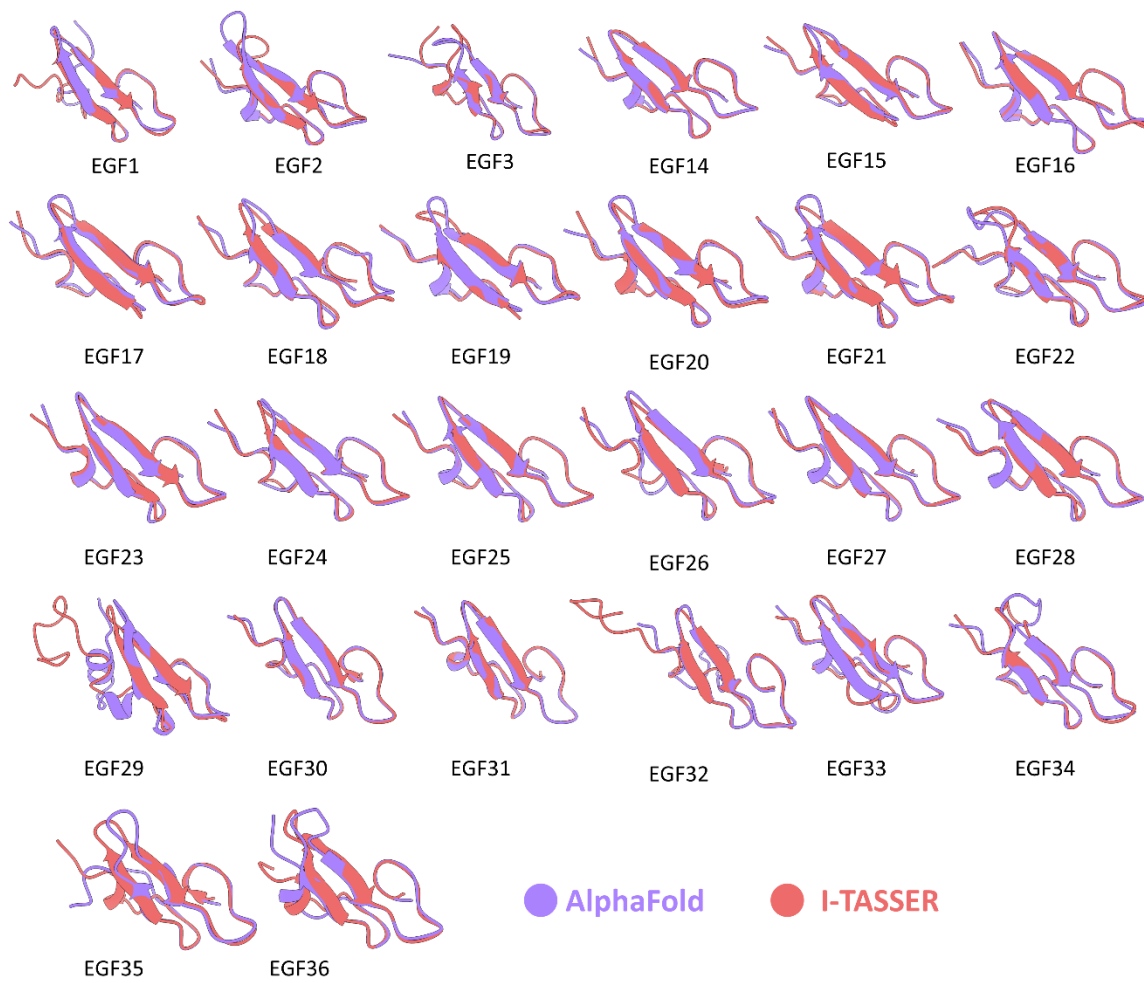
**Figure S4: Comparative assessment of EGF-like repeat patches refinement using OPLS and Charmm36:** The force-field assessment for refinement of EGF-like repeat patches, EGF1-3, 4-13, 14-17, 18-20, 21-24, 25-28, 29-32, and 33-36 using OPLS and Charmm36-AA revealed more defined and stable secondary structure of EGF-like repeats with OPLS.



**Figure S5: AlphaFold predicted models for the human NOTCH1-ECD with unreliable fold:** The top 5 predicted structures by the stand-alone AlphaFold version 2.2.0. are shown in cartoon representation, where the color code refers to the predicted local distance difference test (pLDDT) score. It is a measure of confidence in the correctness of each predicted residue, where pLDDT values of 70 and above are considered to be accurate. The membrane is represented with dotted line around the transmembrane region in order to depict the relative orientation of the NOTCH-ECD. The graph shows low sequence coverage for the predicted models. The color bar represents sequence identity score, where the lowest sequence identity is shown in red. The line plot highlights the pLDDT scores for each residue in the NOTCH-ECD for all 5 models. As shown in the Predicted Aligned Error (PAE) heatmaps, the predicted relative sub-domain orientation of above-mentioned models is observed to be unreliable.

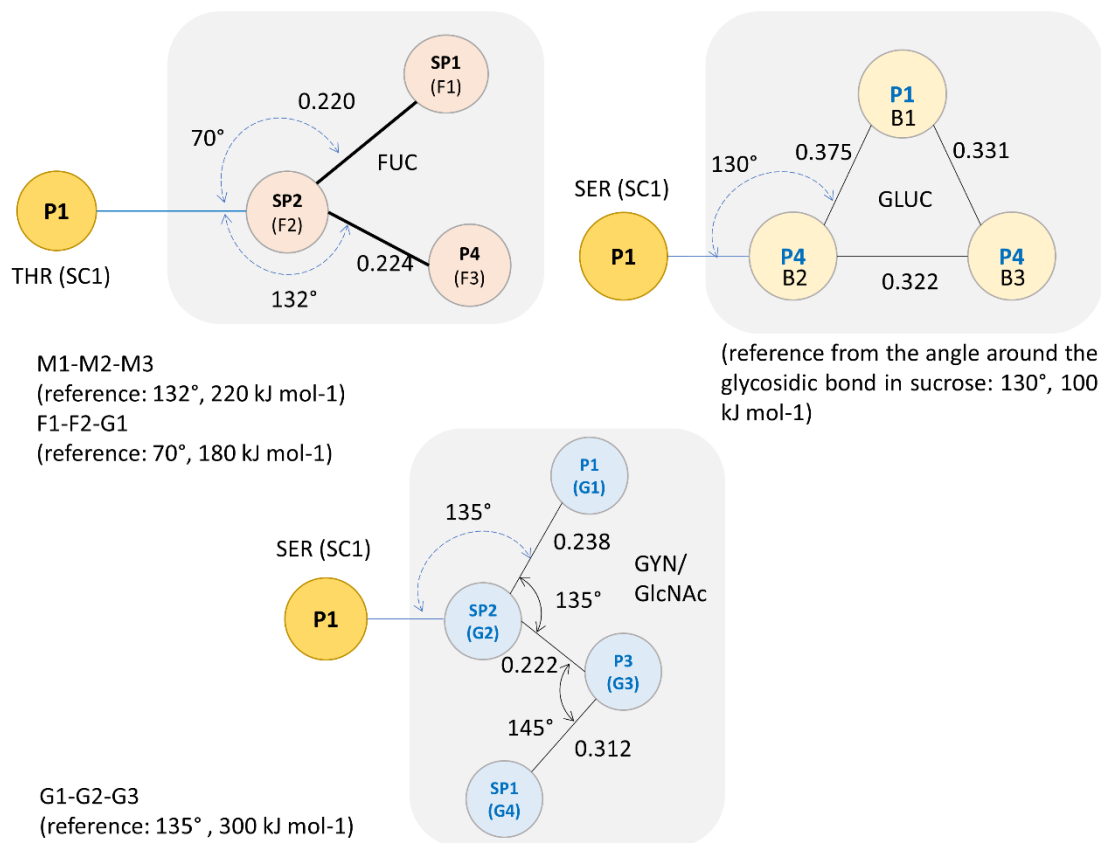


**Figure S6: pLDDT plots for individual EGF-like repeat display high confidence:** The scatter plot for the pLDDT score of each EGF-like repeat model predicted by AlphaFold2 is shown. The blue dots in the scatter plots for individual model of EGF1-3 and EGF14-36 display high pLDDT score.

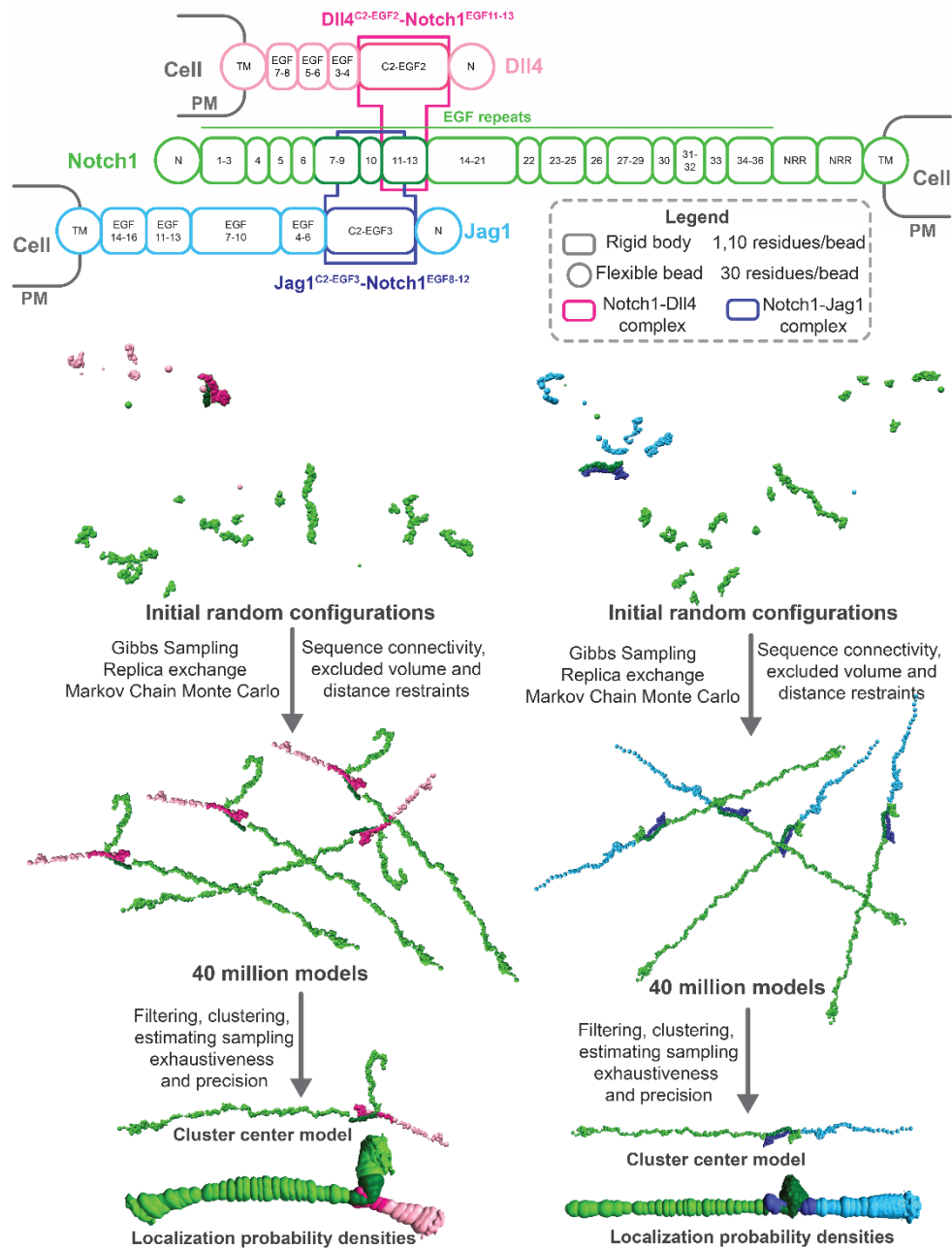


**Figure S7: Comparative assessment of EGF-like repeat models using AlphaFold2 and I-Tasser:** The snapshots showcase high structural alignment between models of EGF-like repeat 1-3 and 14-36 predicted by AlphaFold2 and I-TASSER.

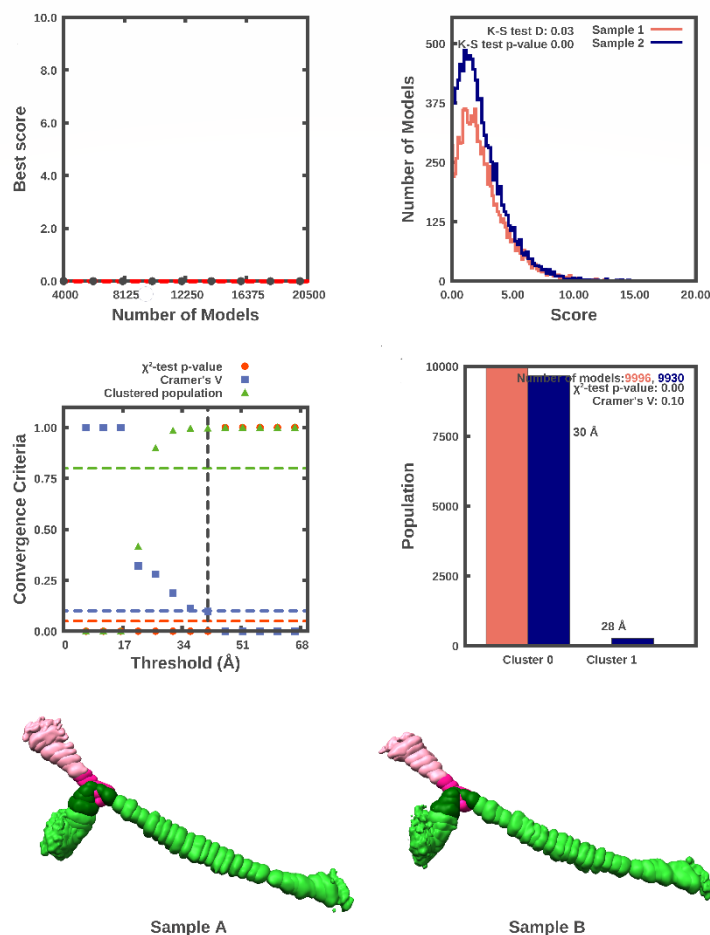




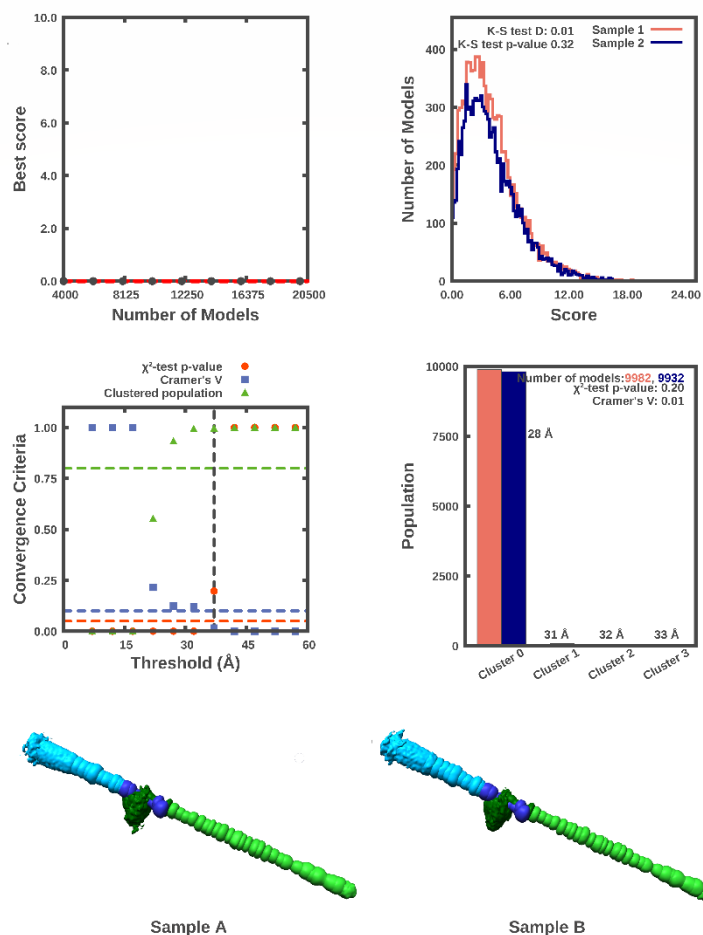
**Figure S8: Glycan parameters for the glucose, fucose and GlcNAc moieties used for glycosylated NOTCH-ECD:** The coarse-grained representation of glucose and fucose is a three-bead representation, whereas GlcNAc is represented as four beads. The respective bond-length, angles, and other parameters have been taken from Shivgan et. al, J. Chem. Inf. Model. 2020.



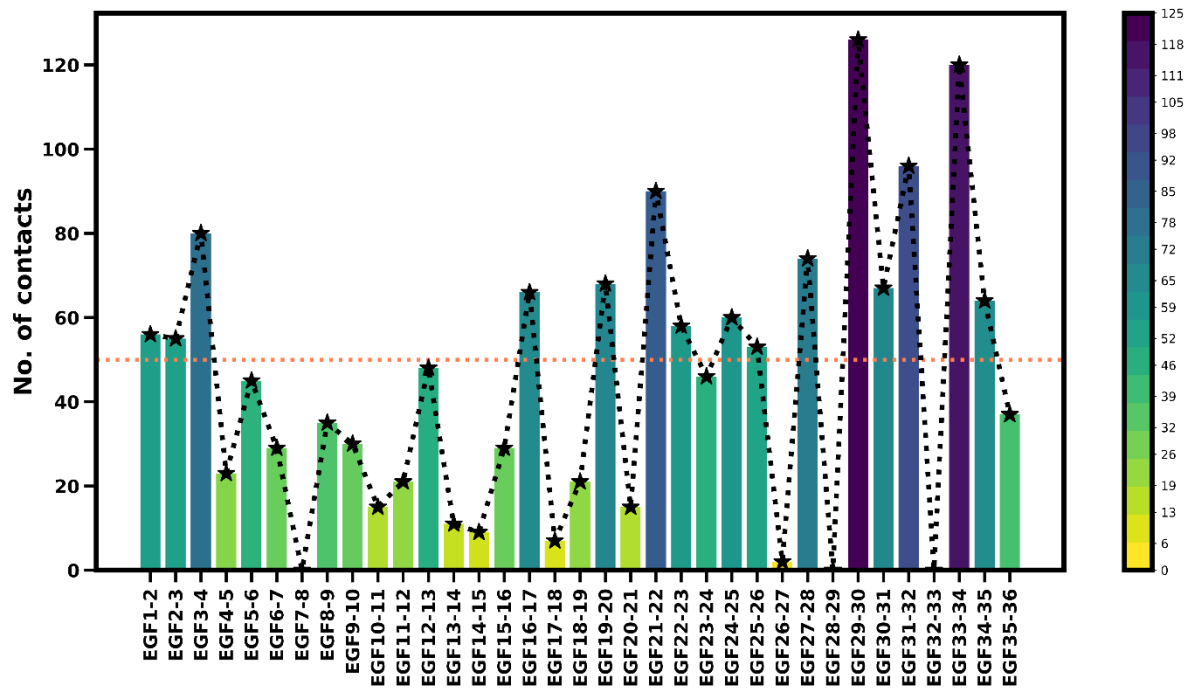
**Figure S9: Integrative structural modelling of NOTCH-DLL4 and NOTCH-JAG1 complexes.** Schematic diagram depicting the structural architecture of NOTCH and its ligands in the integrative model complexes have been shown. For the generation of the coarse-grained models, rigid bodies were defined as 1- and 10- residues per bead and flexible bead regions contained 30-residues per bead. These rigid and flexible bodies are represented by rectangles and circles, respectively. Regions of NOTCH and ligands- DLL4 and JAG1 that are known to form a complex were modelled as rigid bodies and are depicted as dark pink and dark blue boxes, respectively. Workflow for integrative modelling of the NOTCH-DLL4 and NOTCH-JAG1 complex. Starting from initial random configurations, forty million models were sampled using Replica Exchange Markov Chain Monte Carlo algorithm under sequence connectivity, excluded volume and a restraint on the end-to-end length of NOTCH-ECD.



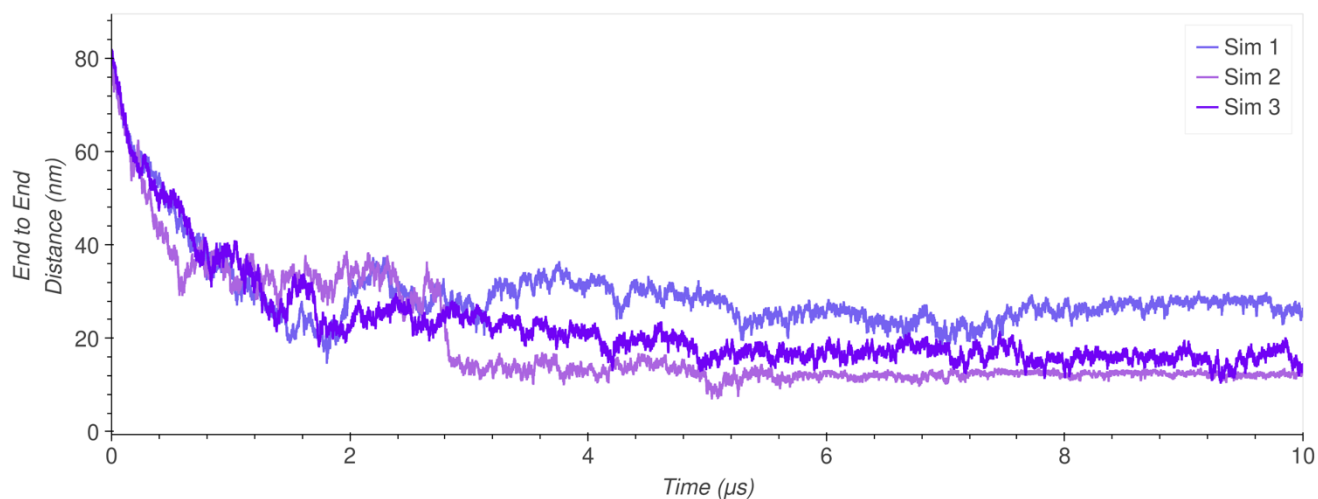
**Figure S10: Sampling exhaustiveness protocol on NOTCH-DLL4 integrative models.** The graph highlights the convergence of the model score for the 20,000 good-scoring models. The scores do not continue to improve as more models are computed essentially independently. The error bar represents the standard deviations of the best scores estimated by iterating sampling of models for 10 cycles. The red dotted line indicates a lower bound reference on the total score. The distribution graph shows the testing similarity of model score between sample 1 (red) and 2 (blue). The difference in the distribution of scores is significant (Kolmogorov-Smirnov two-sample test p-value is  $< 0.05$ ), however the magnitude of the difference is small (Kolmogorov-Smirnov two-sample test statistic  $D$  is  $< 0.3$ ). Hence, the two score distributions are effectively equal. The plot shows three criteria for determining the sampling precision (Y-axis) evaluated as a function of the RMSD clustering threshold (X-axis). The criteria taken are: (i) the p-value is computed using the  $\chi^2$ -test for homogeneity of proportions (red dots), (ii) an effect size for the  $\chi^2$ -test is quantified by the Cramer's  $V$  value (blue squares), and (iii) the population of models in sufficiently large clusters (containing at least 10 models from each sample) is shown as green triangles. The vertical dotted grey line indicates the RMSD clustering threshold at which all three conditions are satisfied (p-value  $> 0.05$ ; dotted red line), Cramer's  $V < 0.10$  (dotted blue line), and the population of clustered models  $> 0.80$  (dotted green line), thus defining the sampling precision of 41 Å. The bar graph depicts the population of models in sample 1 and 2 for the clusters obtained by threshold-based clustering (RMSD threshold = 41 Å). (E-F) The comparison of localization probability densities for NOTCH1-DLL4 models from sample A & B in the major cluster (98% population) are shown. The cross-correlation of the density maps for the two samples is greater than 0.97.



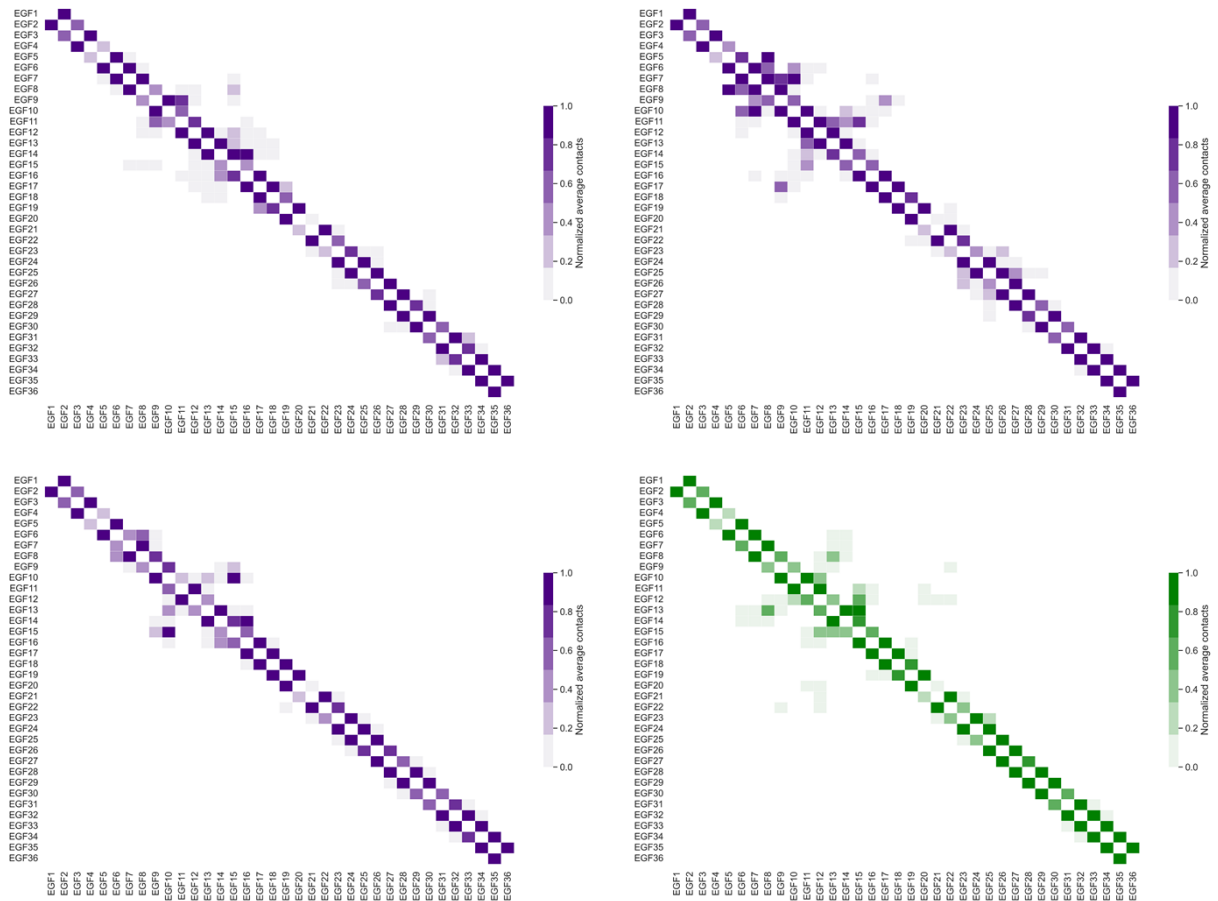
**Figure S11: Sampling exhaustiveness protocol on NOTCH-JAG1 integrative models.** The graph highlights the convergence of the model score for the 20,000 good-scoring models. The scores do not continue to improve as more models are computed essentially independently. The error bar represents the standard deviations of the best scores estimated by iterating sampling of models for 10 cycles. The red dotted line indicates a lower bound reference on the total score. The distribution graph shows the testing similarity of model score between sample 1 (red) and 2 (blue). The difference in the distribution of scores is significant (Kolmogorov-Smirnov two-sample test p-value is  $< 0.05$ ), however the magnitude of the difference is small (Kolmogorov-Smirnov two-sample test statistic  $D$  is  $< 0.3$ ). Hence, the two score distributions are effectively equal. The plot shows three criteria for determining the sampling precision (Y-axis) evaluated as a function of the RMSD clustering threshold (X-axis). The criteria taken are: (i) the p-value is computed using the  $\chi^2$ -test for homogeneity of proportions (red dots), (ii) an effect size for the  $\chi^2$ -test is quantified by the Cramer's  $V$  value (blue squares), and (iii) the population of models in sufficiently large clusters (containing at least 10 models from each sample) is shown as green triangles. The vertical dotted grey line indicates the RMSD clustering threshold at which all three conditions are satisfied (p-value  $> 0.05$ ; dotted red line), Cramer's  $V < 0.10$  (dotted blue line), and the population of clustered models  $> 0.80$  (dotted green line), thus defining the sampling precision of  $37 \text{ \AA}$ . The bar graph depicts the population of models in sample 1 and 2 for the clusters obtained by threshold-based clustering (RMSD threshold =  $41 \text{ \AA}$ ). (E-F) The comparison of localization probability densities for NOTCH-JAG1 models from sample A & B in the major cluster (98% population) are shown. The cross-correlation of the density maps for the two samples is greater than 0.99.



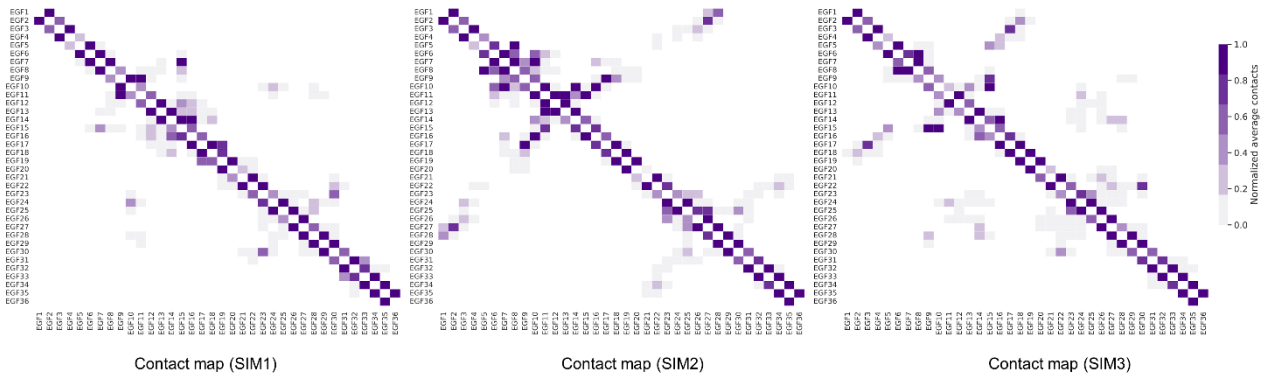
**Figure S12: Lesser intra EGF-like repeat interactions at the NOTCH-ECD N-terminal indicates towards its flexibility.** The bar plot depicts the number of inter-domain contacts between each EGF-like repeat within 5 Å, where the color bar intensity (yellow to blue) represents the increasing number of contacts. A lower number of contacts in the EGF-like repeats proximal to the N-terminus (EGF1-21) are observed indicating the inherently flexible nature of the region.



**Figure S13: Overall dynamics reveal high deviations in the apo NOTCH-ECD.** The line plot highlights the reduction in the end-to-end distance of 36 EGF-like repeats, where the distance was measured between the backbone atoms of ARG20 of EGF1 and the HIS1426 of EGF36 across time, thereby indicating the folding of NOTCH-ECD.

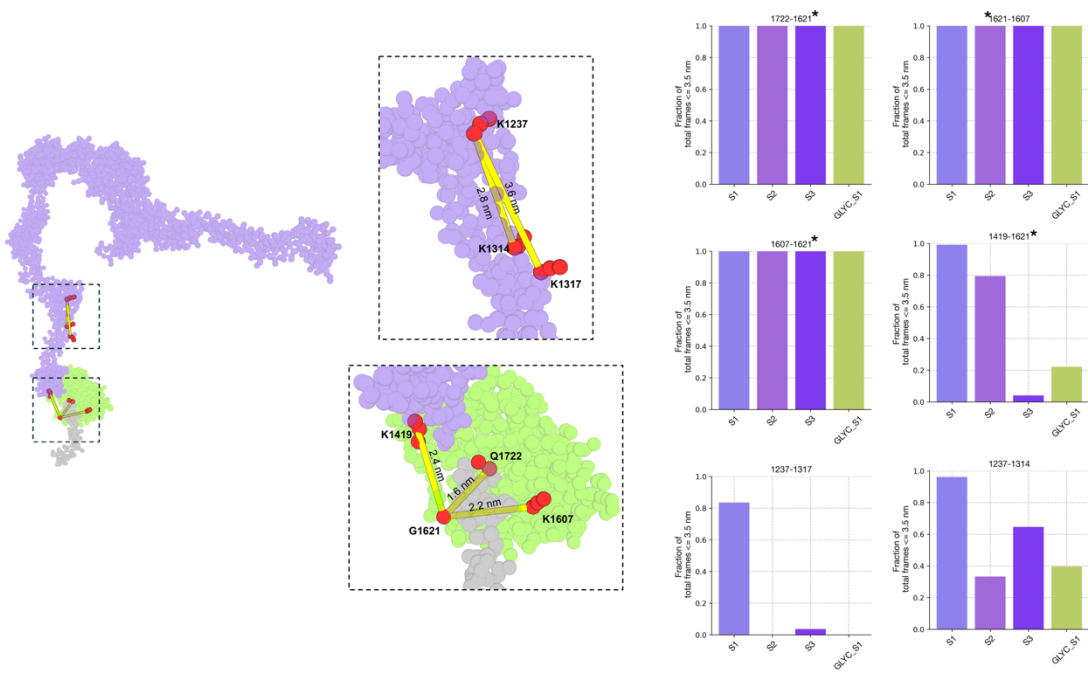


**Figure S14: Distance-based matrix for folded NOTCH in glycosylated vs non-glycosylated form till 1  $\mu$ s.** The distance matrix maps for the three simulations highlighting the contribution of specific EGF-like repeats in the folding of NOTCH-ECD till 1  $\mu$ s have been shown. The color bar (white to dark purple: non-glycosylated NOTCH-ECD and white to dark green: glycosylated NOTCH-ECD) depicts contact occupancy ranging from 0 to 1, where regions making higher contacts in the compact non-glycosylated and glycosylated NOTCH-ECD state can be seen in dark purple and dark green, respectively.

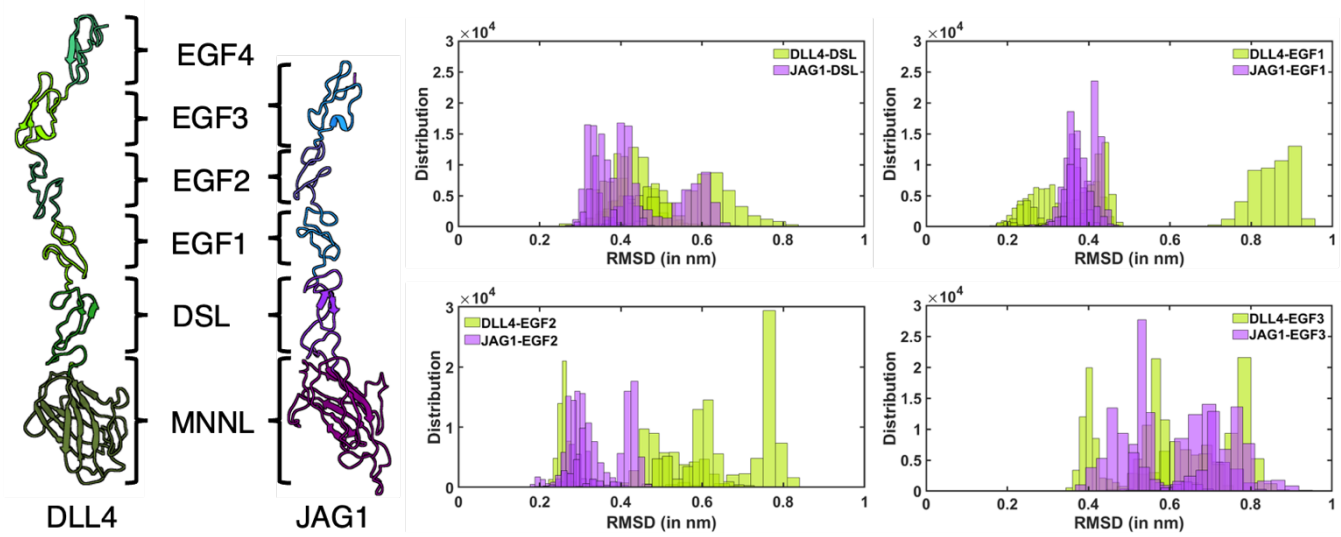


**Figure S15: Distance-based contact matrix reveal specific NOTCH EGF-like repeats regulate compact conformation.** The distance matrix maps for the three simulations highlighting the contribution of specific EGF-like repeats in the folding of NOTCH-ECD have been shown. The color bar (white to dark purple) depicts contact occupancy ranging from 0 to 1, where regions making higher contacts in the compact NOTCH-ECD state can be seen in dark purple.

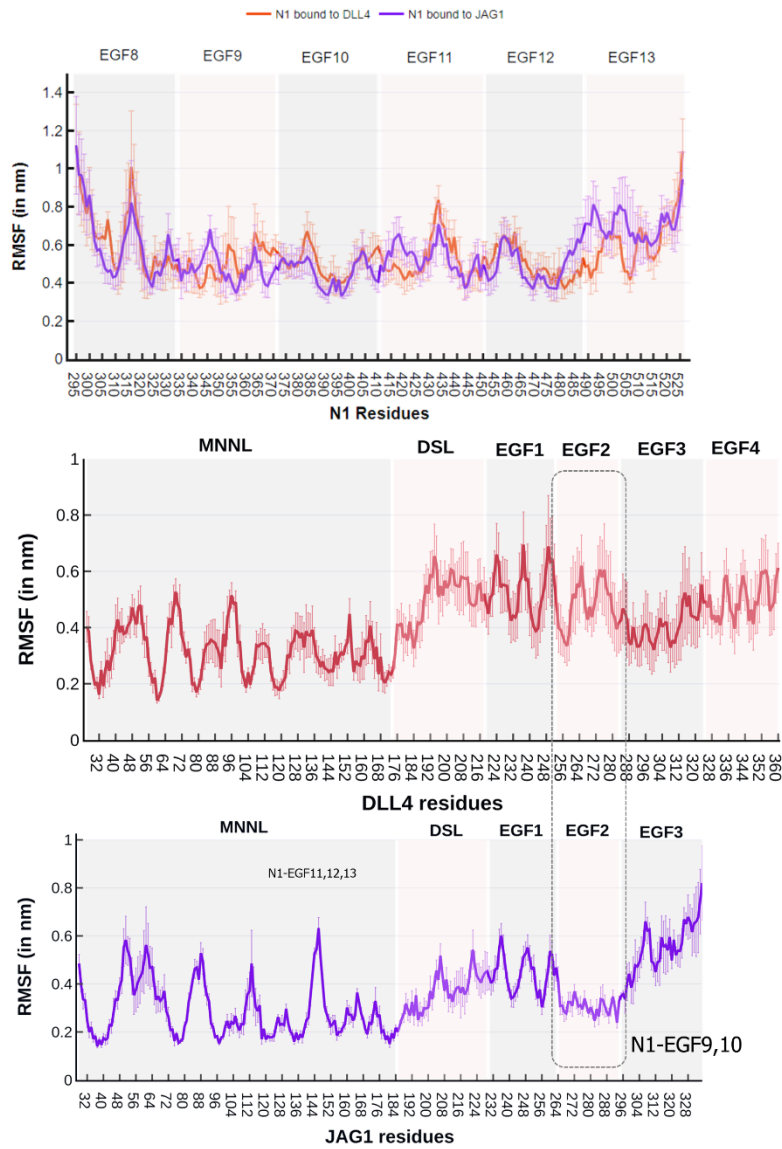




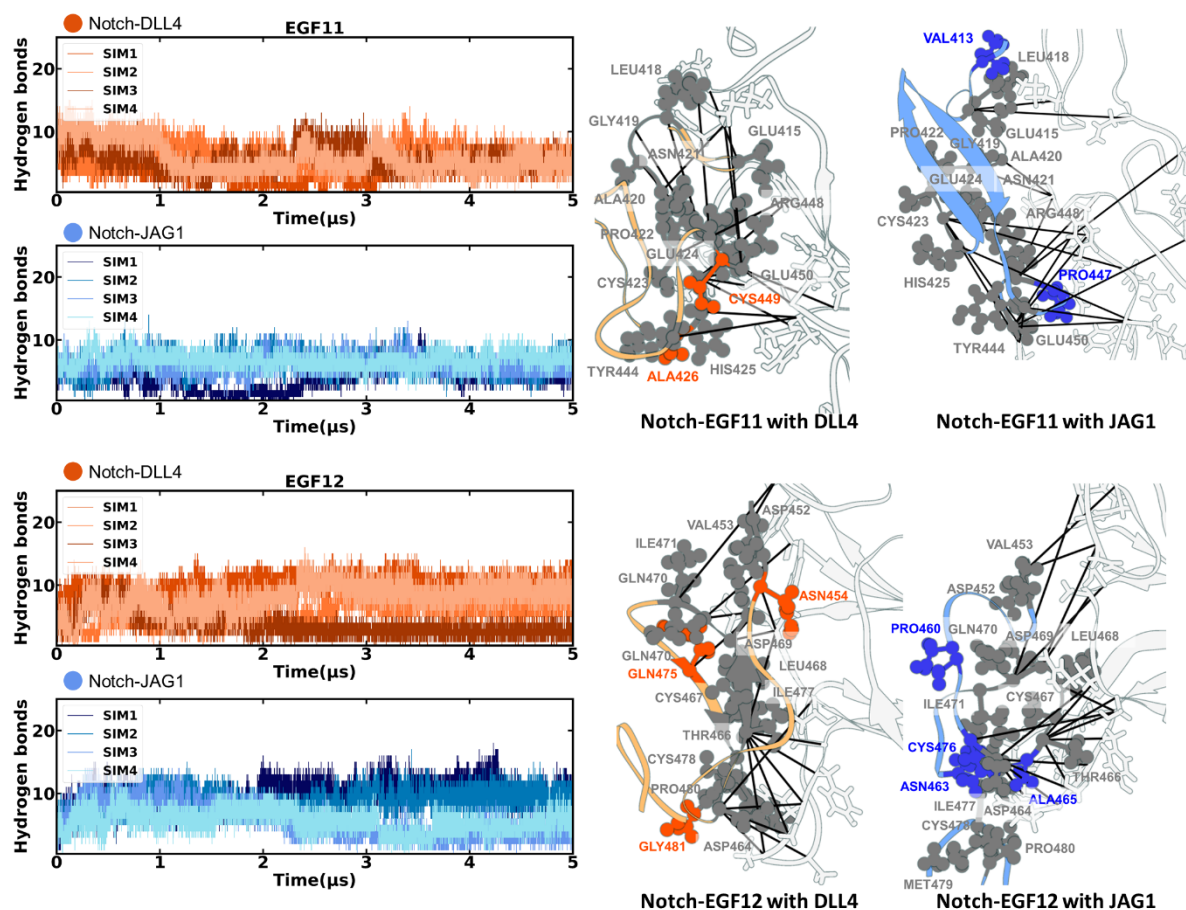
**FigureS16: NOTCH-ECD simulations satisfy intra-NOTCH crosslink data.** The representative snapshot of NOTCH-ECD from MD simulations highlighting the satisfied cross-link pairs is shown. The residue sites are zoomed -in and the corresponding distances are marked. The bar plot highlights the fraction of total frames where two residues have a distance  $\leq 3.5$  nm.



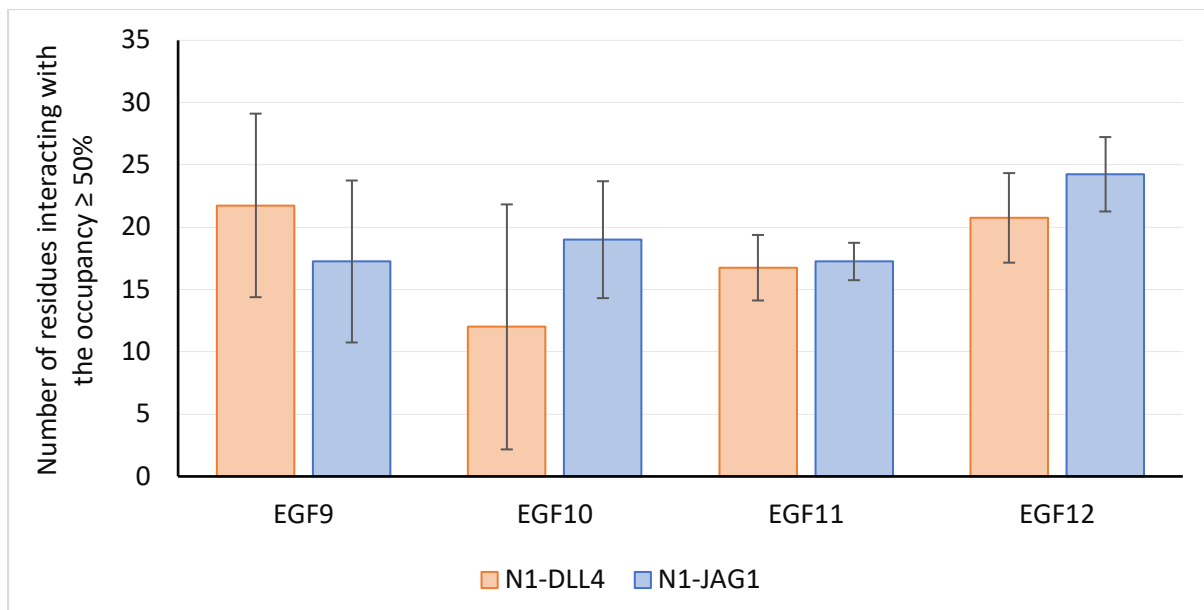
**Figure S17: Structural architecture and dynamics of specific domains in DLL4 and JAG1.** The structural snapshots for DLL4 and JAG1 highlighting their different domains (MNNL, DSL, EGF-like repeats) are shown with gradients of green and blue, respectively. The RMSD for individual domains of DLL4 (in green) and JAG1 (in purple) are shown via frequency distribution plots.



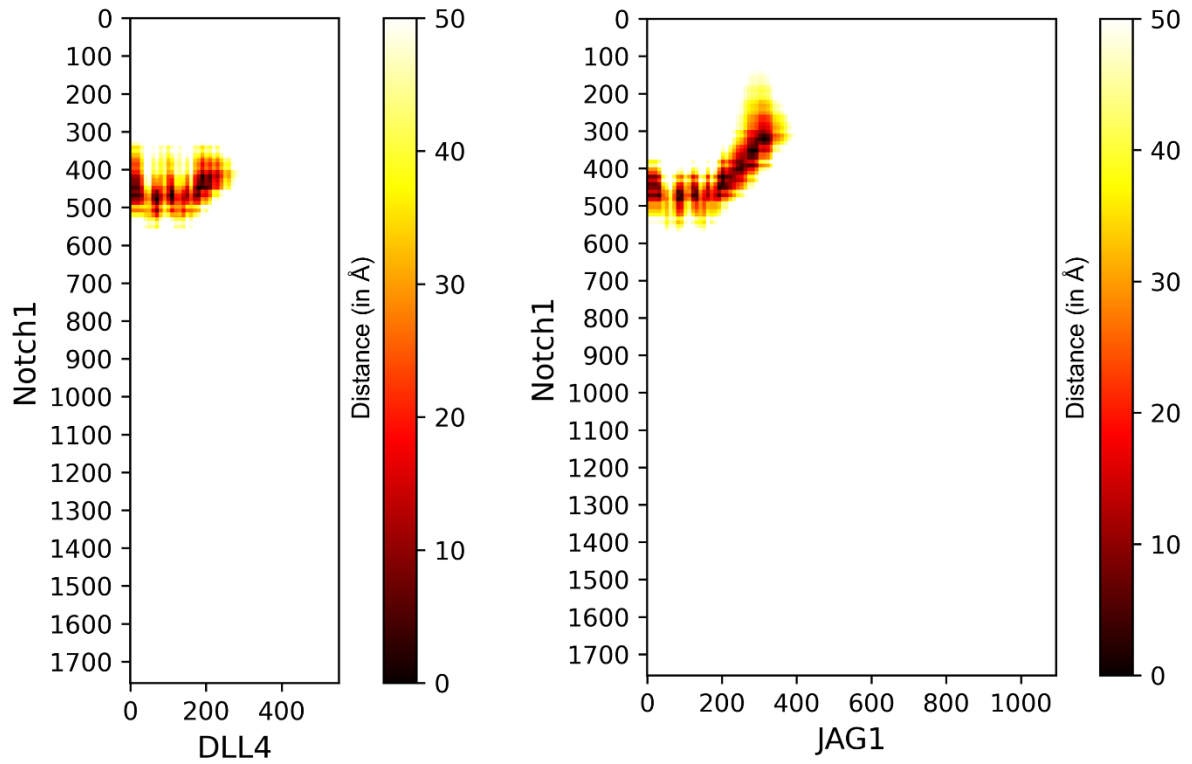
**Figure S18: Domain-wise fluctuations in NOTCH, DLL4 and JAG1 across the four 5 μs simulations.** The line plot shows region-wise root mean square fluctuations of NOTCH, DLL4 and JAG1, where residue-wise fluctuation across 4 replicates have been shown as standard error bars.



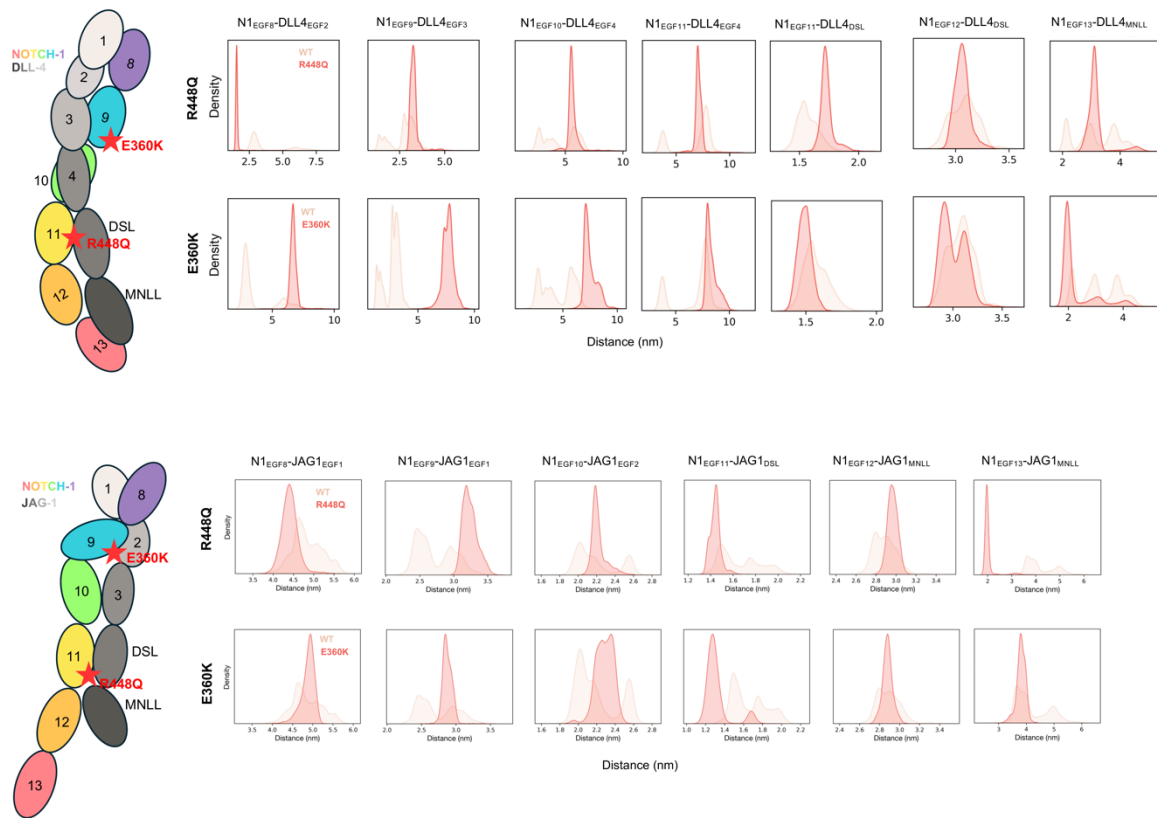
**Figure S19: NOTCH EGF11 and 12 emerge as common contributor in NOTCH-ligand complexes.** The average number of hydrogen bonds mediated by NOTCH EGF11 and 12 with DLL4 and JAG1 across all 4 simulations have been shown with the line graph, where NOTCH EGF11 and 12 show similar number of H-bonds with DLL4 and JAG1. The zoomed-in structural snapshots for NOTCH EGF-like repeats 11 and 12 in NOTCH-DLL4 and NOTCH-JAG1 complexes are shown, where contributing residues (occupancy  $\geq 50\%$ ) are highlighted in sphere representation. Amongst these, the common residues of NOTCH shared between DLL4 and JAG1 are highlighted in grey, while those exclusively found in NOTCH-DLL4 and NOTCH-JAG1 are shown in dark orange and dark blue, respectively. The hydrogen bond pairs reported in the crystal structures of NOTCH-DLL4 and NOTCH-JAG1 were also observed with low occupancies across all runs in addition to several new contacts mediated by the same residues. A few key interactions occurring with high occupancies were observed where GLU424 and GLU450 of NOTCH EGF11 and ASP464 of NOTCH-EGF12 actively mediate hydrogen bonds with both DLL4 and JAG1. GLU424 mediates interaction with DSL residues- LYS189, ARG191, and TYR198 in DLL4 and ARG201 and ARG 203 of JAG1. Similarly, GLU450 interacts with ARG186 and LYS218 in DLL4 and JAG1, respectively. Four key interactions between NOTCH-EGF12 and MNNL of DLL4 were observed as 452ASP-179TYR, ASP464-GLN66, ASP464-THR68 and ASP469-THR110. While ASP464 of N1 EGF12 forms contact with ARG85 in JAG1.



**Figure S20: Varied participation of NOTCH EGF-like repeats in NOTCH-DLL4 and NOTCH-JAG1 complexes.** The bar graph represents the number of NOTCH residues from EGF9-12 with an occupancy of  $\geq 50\%$  participating in contact formation with DLL4 and JAG1 across all simulations.



**Figure S21: Distance matrix plot displaying the interacting regions of full-length NOTCH-ECD with ECD of DLL4 and JAG1.** The matrix plot shows the interacting regions in the integrative NOTCH-DLL4 and NOTCH-JAG1 complexes, where an extended region of NOTCH EGF5-7 can be seen in proximity when bound to JAG1.



**Figure S22: Ligand binding mutations introduce spatial reorientation of EGF-like repeats.** The NOTCH and ligand (DLL4 and JAG1) EGF pairwise distances are plotted as probability distribution plots. The distances are compared between WT and the mutants (R448Q, E360K).

**Table S1.** The respective templates for individual EGF-like repeats and the C-score for the best predicted structure by I-TASSER:

<b>EGF-like repeat with unknown structure</b>	<b>Templates used for model construction</b>	<b>C-score for the predicted EGF-like repeat (ranges between -5 to 2)</b>
EGF1	EGF6	0.6
EGF2	EGF4	0.63
EGF3	EGF7	1.19
EGF14	EGF12	1.52
EGF15	EGF12	1.49
EGF16	EGF5	1.51
EGF17	EGF13	1.47
EGF18	EGF7	1.39
EGF19	EGF12	1.55
EGF20	EGF5	1.52
EGF21	EGF12	1.56
EGF22	EGF6	1.15
EGF23	EGF5	1.53
EGF24	EGF7	1.49
EGF25	EGF12	1.55
EGF26	EGF7	1.42
EGF27	EGF5	1.56
EGF28	EGF8	1.32
EGF29	EGF6	0.39
EGF30	EGF8	1.44
EGF31	EGF12	1.56
EGF32	EGF7	0.9
EGF33	EGF5	1.24
EGF34	EGF12	1.14



EGF35	EGF4	1.37
EGF36	EGF4	0.78

**Table S2.** The quality check scores from SAVES v5.0 for each modelled EGF-like repeat:

<b>EGF MODELS</b>	<b>VERIFY 3D*</b>	<b>ERRAT**</b>
EGF1	100%	100
EGF2	85.37%	100
EGF3	71.05%	83.33
EGF14	100%	100
EGF15	100%	100
EGF16	100%	100
EGF17	100%	100
EGF18	100%	100
EGF19	100%	100
EGF20	100%	100
EGF21	70.27%	100
EGF22	100%	92
EGF23	100%	100
EGF24	67.57%	100
EGF25	100%	100
EGF26	100%	96.29
EGF27	100%	100
EGF28	100%	92.30
EGF29	57.45%	96.96
EGF30	100%	100
EGF31	100%	100
EGF32	75.56%	100
EGF33	100%	93.54
EGF34	97.50%	88
EGF35	62.16%	100
EGF36	65%	80

\* At least 80% of the amino acids have scored  $\geq 0.2$  in the 3D/1D profile.

\*\* Overall Quality Factor

**Table S3.** Type and site of glycosylation present in NOTCH-ECD

S. No.	EGF repeat	Glycan type	Residue no.	Residue name
1	EGF2	O-glucose	65	Serine
2	EGF2	O-fucose	73	Threonine
3	EGF3	O-fucose	116	Threonine
4	EGF4	O-glucose	146	Serine
5	EGF5	O-fucose	194	Threonine
6	EGF6	O-fucose	232	Threonine
7	EGF8	O-fucose	311	Threonine
8	EGF9	O-GlcNAc	341	Serine
9	EGF9	O-fucose	349	Threonine
10	EGF10	O-glucose	378	Serine
11	EGF11	O-glucose	435	Serine
12	EGF12	O-glucose	458	Serine
13	EGF12	O-fucose	466	Threonine
14	EGF13	O-glucose	496	Serine
15	EGF14	O-glucose	534	Serine
16	EGF16	O-glucose	609	Serine
17	EGF16	O-fucose	617	Threonine
18	EGF17	O-glucose	647	Serine
19	EGF18	O-fucose	692	Threonine
20	EGF19	O-glucose	722	Serine
21	EGF20	O-glucose	759	Serine
22	EGF20	O-fucose	767	Threonine
23	EGF20	O-GlcNAc	784	Serine
24	EGF21	O-glucose	797	Serine
25	EGF21	O-fucose	805	Threonine
26	EGF25	O-glucose	951	Serine
27	EGF26	O-fucose	997	Threonine
28	EGF27	O-glucose	1027	Serine
29	EGF28	O-fucose	1035	Threonine
30	EGF28	O-glucose	1065	Serine
31	EGF30	O-fucose	1159	Threonine
32	EGF31	O-glucose	1189	Serine
33	EGF31	O-fucose	1197	Threonine
34	EGF33	O-glucose	1273	Serine
35	EGF35	O-fucose	1362	Threonine
36	EGF36	O-fucose	1402	Threonine

**Table S4.** Representation and degrees of freedom for integrative modeling of the (A) NOTCH1-DLL4 (B) NOTCH1-Jag1 complex

**(A) NOTCH-DLL4**

<b>Protein</b>	<b>Residue range</b>	<b>Known atomic structure (PDB or homology template PDB; modeled at 1 &amp; 10 residues per bead) or unknown structure (designated as Bead; modeled at 30 residues per bead)</b>	<b>Domains comprising a single rigid body</b>	<b>Residues in rigid body</b>	
NOTCH	1-19	Bead	-	-	
	20-1735	NOTCH-ECD structure (this study)	EGF1-3	20-139	
			EGF4	140-177	
			EGF5	178-217	
			EGF6	218-256	
			EGF7-9	257-371	
			EGF10	372-411	
			Homology, PDB ID: 4XLW	EGF11-13	412-526
			NOTCH-ECD structure (this study)	EGF14-21	527-828
				EGF22	829-868
				EGF23-25	869-982
		EGF26		983-1020	
		EGF27-29		1021-1144	
		EGF30	1145-1182		
		EGF31-32	1183-1266		

			EGF33	1267-1306
			EGF34-36	1307-1448
			NRR_1	1449-1592
			NRR_2	1593-1735
	1736-1756	Bead	-	-
DLL4	1-25	Bead	-	-
	26-285	Homology, PDB ID: 4XLW	C2-EGF2	26-285
	286-306	Bead	-	-
	307-443	Homology, PDB ID: 4CC0	EGF3-4	307-361
			EGF5-6	362-443
	444-529	Homology, PDB ID: 6M3B	EGF7-8	444-529
530-550	Bead	-	-	

### (B) NOTCH-JAG1

Protein	Residue range	Known atomic structure (PDB or homology template PDB; modeled at 1 & 10 residues per bead) or unknown structure (designated as Bead; modeled at 30 residues per bead)	Domains comprising a single rigid body	Residues in rigid body
NOTCH	1-19	Bead	-	-
	20-1735	NOTCH-ECD structure (this study)	EGF1-3	20-139
			EGF4	140-177
			EGF5	178-217

			EGF6	218-256
			EGF7	257-294
		Homology, PDB ID: 5UK5	EGF8-12	295-489
		NOTCH-ECD structure (this study)	EGF13-21	490-828
			EGF22	829-868
			EGF23-25	869-982
			EGF26	983-1020
			EGF27-29	1021-1144
			EGF30	1145-1182
			EGF31-32	1183-1266
			EGF33	1267-1306
			EGF34-36	1307-1448
			NRR_1	1449-1592
			NRR_2	1593-1735
	1736-1756		Bead	-
JAG1	1-27	Bead	-	-
	28-335	Homology, PDB ID: 5UK5	C2-EGF3	28-335
	336-628	Homology, PDB ID: 1N7D	EGF4-6	336-449
			EGF7-10	450-629
	629-856	Homology, PDB ID: 6POG	EGF11-13	630-742

			EGF14-16	743-845
	846-1093	Bead	-	-

**Table S5.** Intra-NOTCH crosslinks and distance between respective residue pairs across NOTCH-ECD simulations

Position A in Mouse	Corresponding position A in Human	Position B in Mouse	Corresponding position B in Human	S1 (nm) Mean+/- SD	S2 (nm)	S3 (nm)	Glyc Sim (nm)	Remarks
K1419	K1419	K1314	K1314	5.62+/- 0.18	5.63+/- 0.17	5.54+/- 0.22	5.61+/- 0.20	
K1237	K1237	K1314	K1314	2.76+/- 0.30	3.56+/- 0.17	3.36+/- 0.31	3.53+/- 0.21	Satisfied
K1628	K1628 Not in X-ray structure (Picking the nearest available residue G1621)	K1607	K1607	2.09+/- 0.07	2.10+/- 0.07	2.09+/- 0.07	2.15+/- 0.07	Satisfied
K1607	K1607	K1632	K1632 Not in X-ray structure (Picking the nearest available residue G1621)	2.09+/- 0.07	2.10+/- 0.07	2.09+/- 0.07	2.15+/- 0.07	Satisfied
K1237	K1237	K1317	K1317	3.27+/- 0.35	4.23+/- 0.19	4.05+/- 0.31	4.23+/- 0.21	Satisfied in S1 (Others are closer to 4nm)
K1632	1632 Not in X-ray structure (Picking the nearest available residue G1621)	K1314	K1314	7.57+/- 0.39	7.05+/- 0.69	6.28+/- 0.79	7.62+/- 0.55	

K1498	K1498	K1314	K1314	6.09+/- 0.67	6.16+/- 0.38	5.97+/- 0.43	4.94+/- 1.20	Decrease in Glycan Sim
K1712	Q1722	K1628	K1628 Not in X-ray structure (Picking the nearest available residue G1621)	1.59+/- 0.09	1.61+/- 0.08	1.59+/- 0.08	1.64+/- 0.09	Satisfied
K1107	K1107	K1314	K1314	7.65+/- 1.08	9.29+/- 0.37	8.97+/- 0.97	9.95+/- 0.51	
K1419	K1419	K1628	K1628 Not in X-ray structure (Picking the nearest available residue G1621)	2.39+/- 0.41	3.54+/- 0.48	4.17+/- 0.34	3.67+/- 0.32	Satisfied