

CASP 13

Predicting Contacts



Albert Einstein College of Medicine

Assessor: András Fiser

Department of Systems and Computational Biology

Department of Biochemistry

Possible questions

- Does contact prediction accuracy correlate with that of structure modeling?
- How well you did among yourselves?
- How well you did compared to previous CASPs?
- Some insight analysis:
 - Are you capturing the same set of contacts?
 - Are there particular types of contacts that you are getting accurately?
 - How important is the quality of sequence information?

Best structure prediction (out of 98) vs. Best contact predictions (out of 46)

FM and TBM/FM	FM	Contacts	Contacts only
G043	G043	X	G498 (6)
G322	G322	X	G032 *
G089	G089	G089 (20)	G180 *
G145	G145	X	G323 *
G224	G224	G224 (11)	G491 *
G261	G498	G498 (1)	G106 *
G354	G261	X	G164 (46)
G498	G354	X	G189 *
G197	G197	X	G352 *
G460	G324	X	G125 *
G324	G196	X	G224 (5)
G135	G208	X	G036 *
G196	G460	X	G392 *
G055	G135	X	G351 (54)
G418	G055	X	G122 (67)
G117	G117	X	G386 *
G208	G418	X	G475 *
G274	G366	X	G154 *
G086	G192	X	G292 *
G192	G274	X	G089 (3)
G071	G086	X	G430 *
G222	G457	X	G041 (63)
G044	G044	X	G091 *

Best structure prediction (out of 98) vs. Best contact predictions (out of 46)

FM and TBM/FM	FM	Contacts	Contacts only
G043	G043	X	G498 (6)
G322	G322	X(G036)(12)	G032 *(2)G322,
G089	G089	G089 (20)	G180 *(2)G322
G145	G145	X(G032)(2)	G323 *(2)G322
G224	G224	G224 (11)	G491 *(16)G117
G261	G498	G498 (1)	G106 *
G354	G261	X(G180, G32)(3,2)	G164 (46)
G498	G354	X(G229)(39)	G189 *
G197	G197	X	G352 *
G460	G324	X(G498)	G125 *
G324	G196	X	G224 (5)
G135	G208	X	G036 *(2) or (50) _{G116}
G196	G460	X	G392 *
G055	G135	X	G351 (54)
G418	G055	X	G122 (67)
G117	G117	X(G491)	G386 *
G208	G418	X	G475 *
G274	G366	X	G154 *
G086	G192	X	G292 *
G192	G274	X	G089 (3)
G071	G086	X	G430 *
G222	G457	X	G041 (63)
G044	G044	X	G091 *

Best structure prediction (out of 98) vs. Best contact predictions (out of 46)

FM and TBM/FM	FM	Contacts	Contacts only
G043	G043	X	G498 (6)
G322	G322	X	G032 *
G089	G089	G089 (20)	G180 *

Difficult to establish clear relation between contact and structure prediction
-> we do not know how well one could perform with a "top" contact prediction

G354	G261	X	G164 (46)
G408	G254	X	G180 *

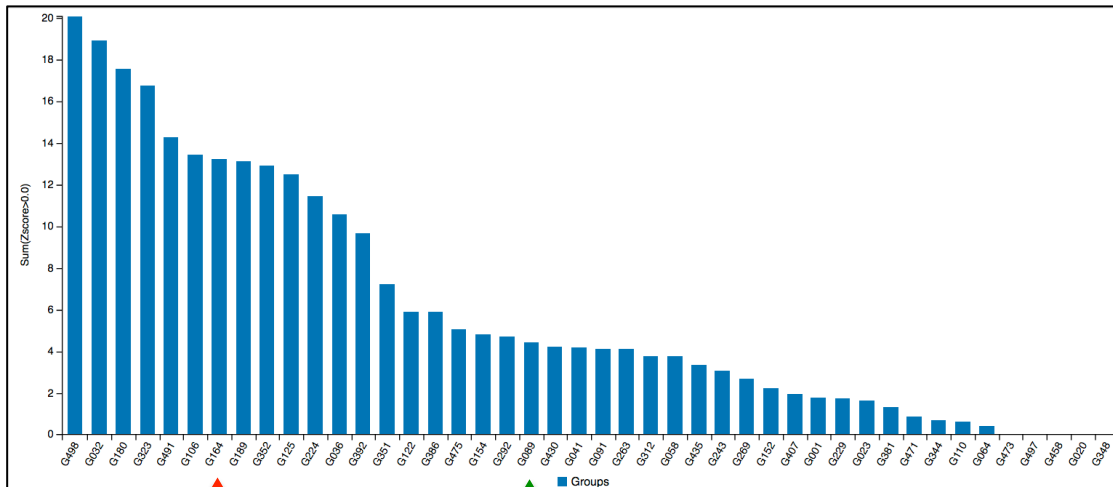
Some top performing structure prediction groups did not submit contact prediction
-> we do not know if they have a better contact prediction than others

G324	G190	X	G224 (5)
G135	G208	X	G036 *
G196	G460	X	G392 *
G055	G135	X	G351 (54)
G418	G055	X	G122 (67)
G117	G117	X	G386 *
G208	G418	X	G475 *
G274	G366	X	G154 *
G086	G192	X	G292 *
G192	G274	X	G089 (3)
G071	G086	X	G430 *
G222	G457	X	G041 (63)
G044	G044	X	G091 *

Best structure prediction (out of 98) vs. Best contact predictions (out of 46)

Among groups that have submitted both structure and contact prediction: Surprising inconsistencies!!

It is important to know how to use contact information!
And/Or
Contact information is not as important as one thought

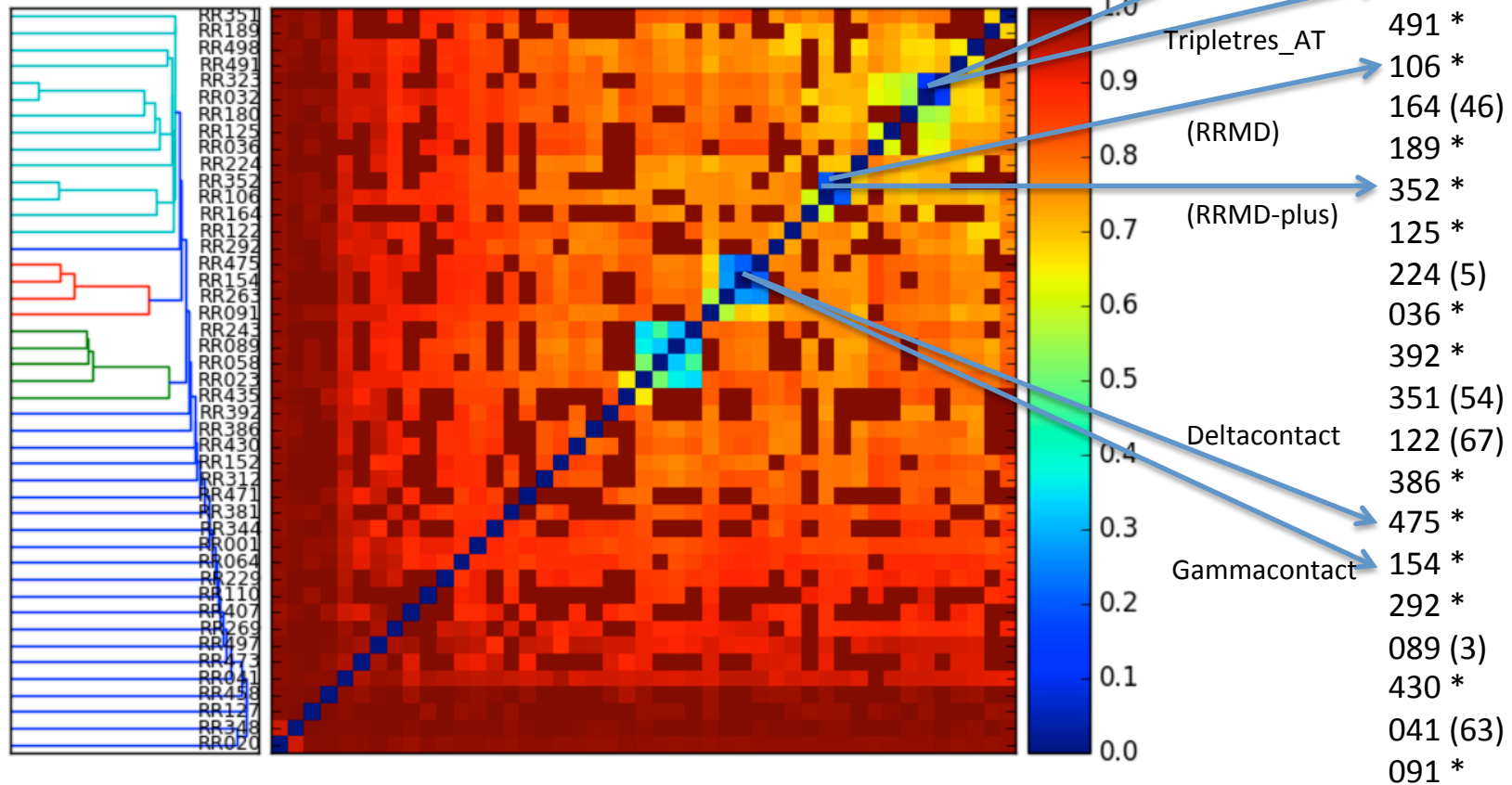


Contacts	Contacts only
X	498 (6)
X	032 *
089 (20)	180 *
X	323 *
224 (11)	491 *
498 (1)	106 *
X	164 (46)
X	189 *
X	352 *
X	125 *
X	224 (5)
X	036 *
X	392 *
X	351 (54)
X	122 (67)
X	386 *
X	475 *
X	154 *
X	292 *
X	089 (3) ^{89 submitted: 30/31 targets...}
X	430 *
X	041 (63)
X	091 *

Are we predicting different contacts? Jaccard distance (“1-Intersection over Union”)

$$d_j = \frac{|A \cup B| - |A \cap B|}{|A \cup B|}$$

0 (same) < d_j < 1 (different)



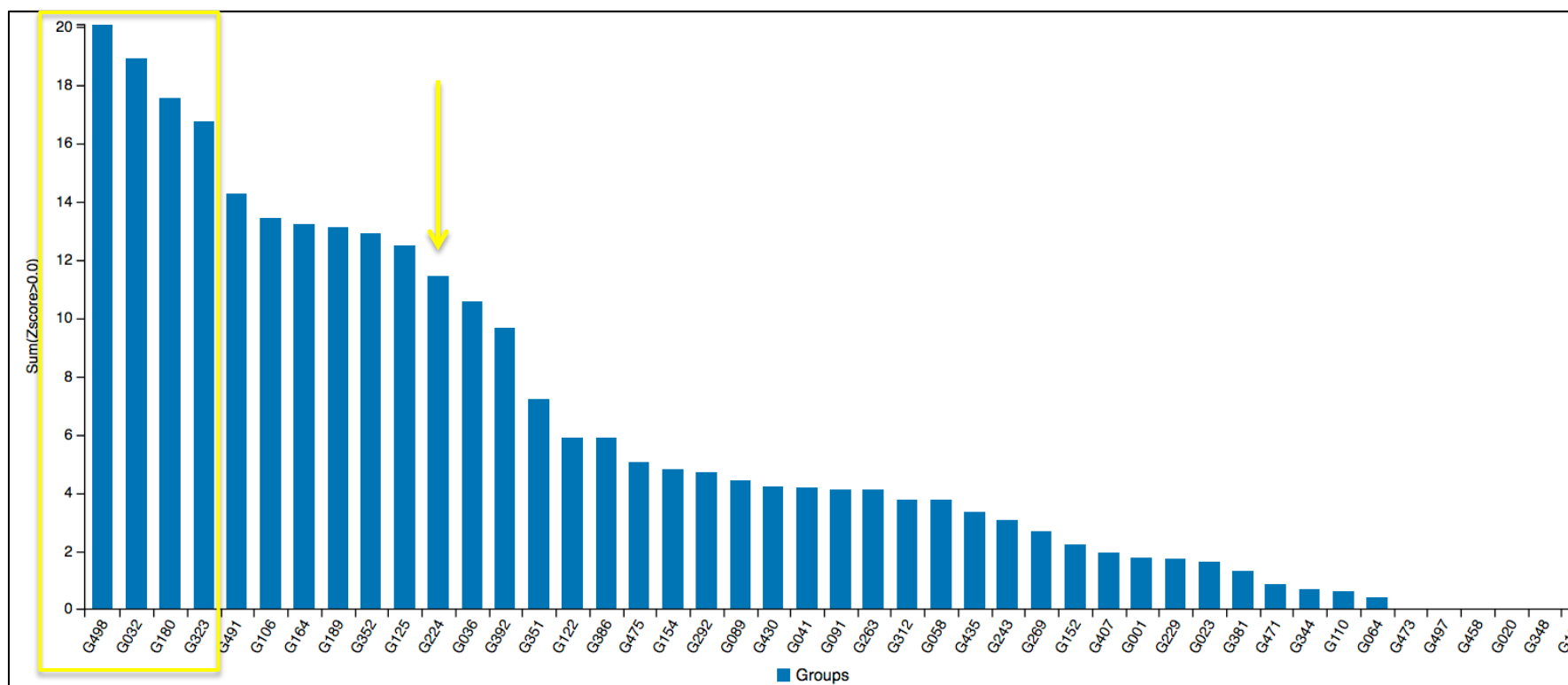
Top L/5 number of contacts, L is the length of sequence

Performance using different criteria

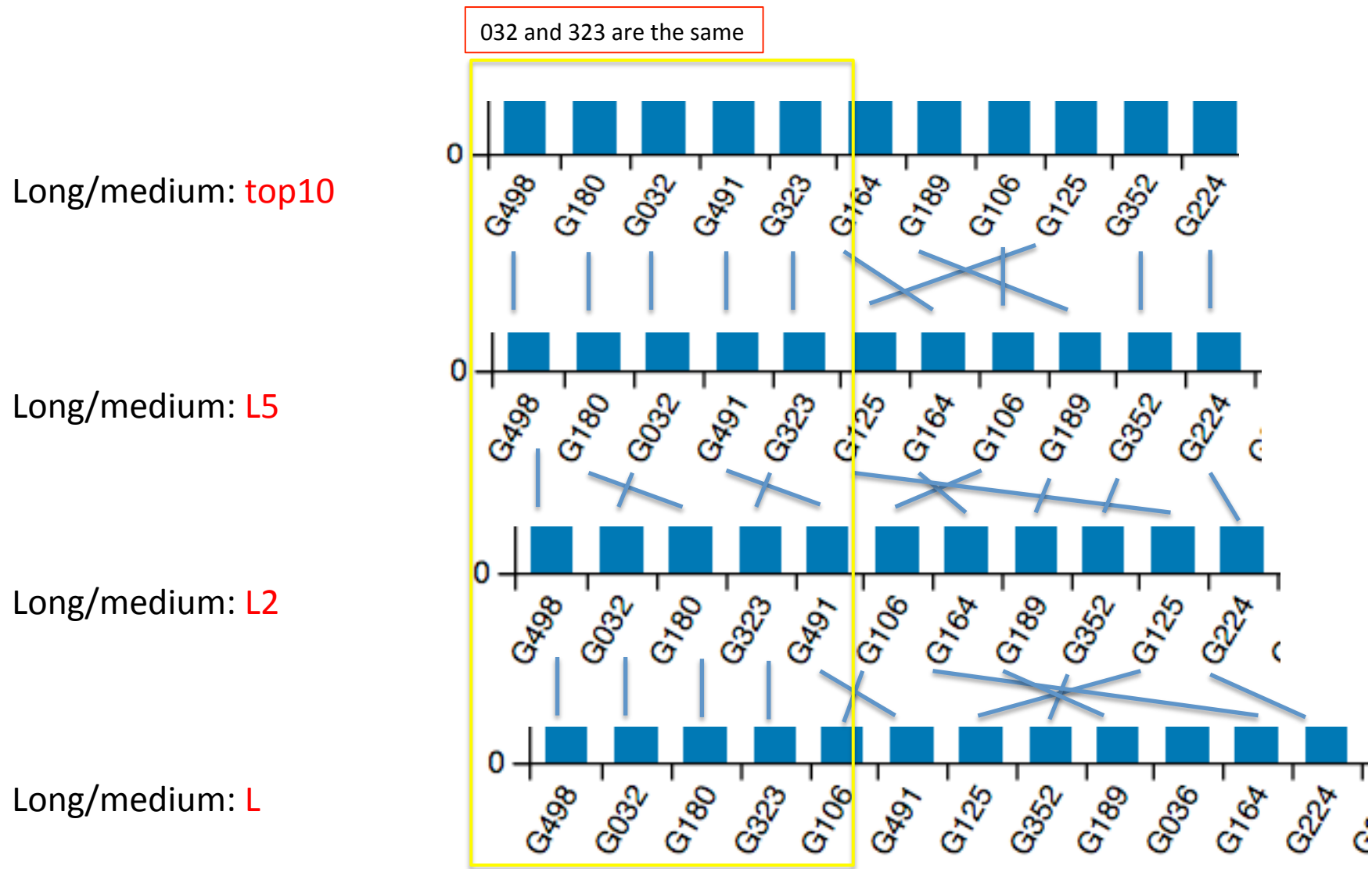
Models (FM or FM+TBM) **X** Contacts (top10 or L/5 or L/2 or L or FL) **X** probability (0 or 0.5) **X** contact definition (medium/long; long; extra long) => 60 combinations

evaluated by either: using F1; Precision/Recall; Z-score sum or Z-score average etc.

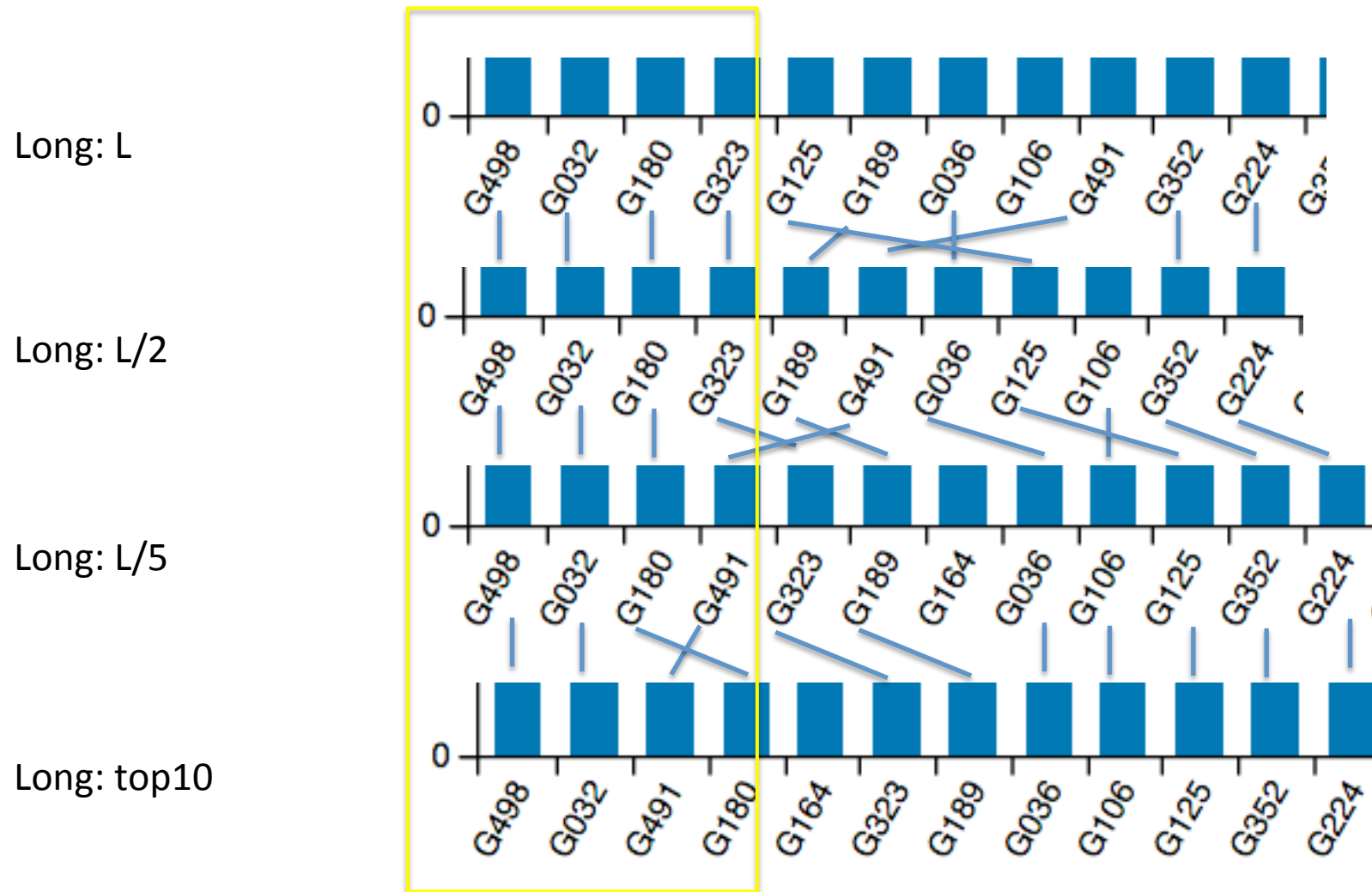
Long/medium contacts, FM only, Zscore >0



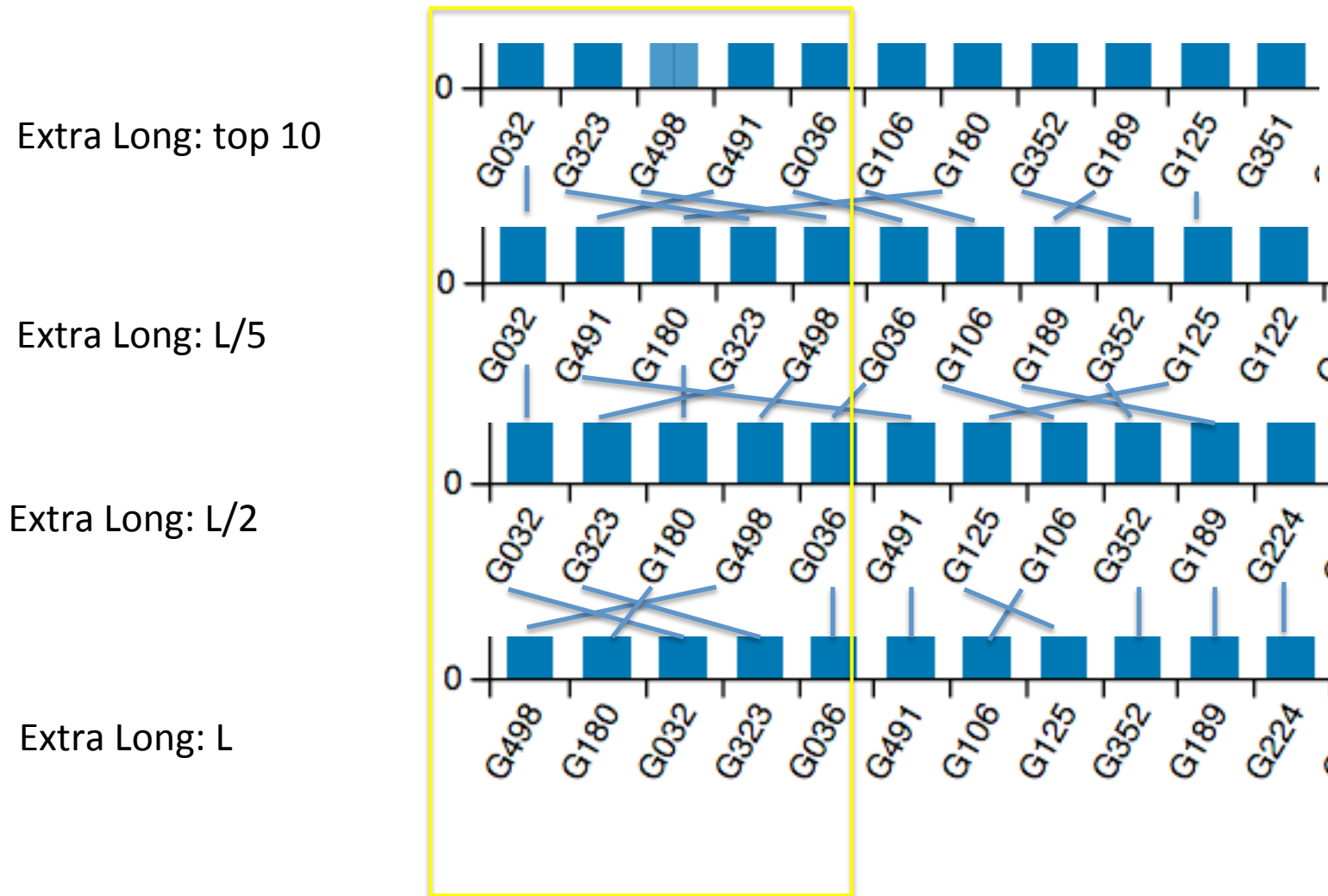
Long/medium contacts (FM only), sum Zscore (>0)



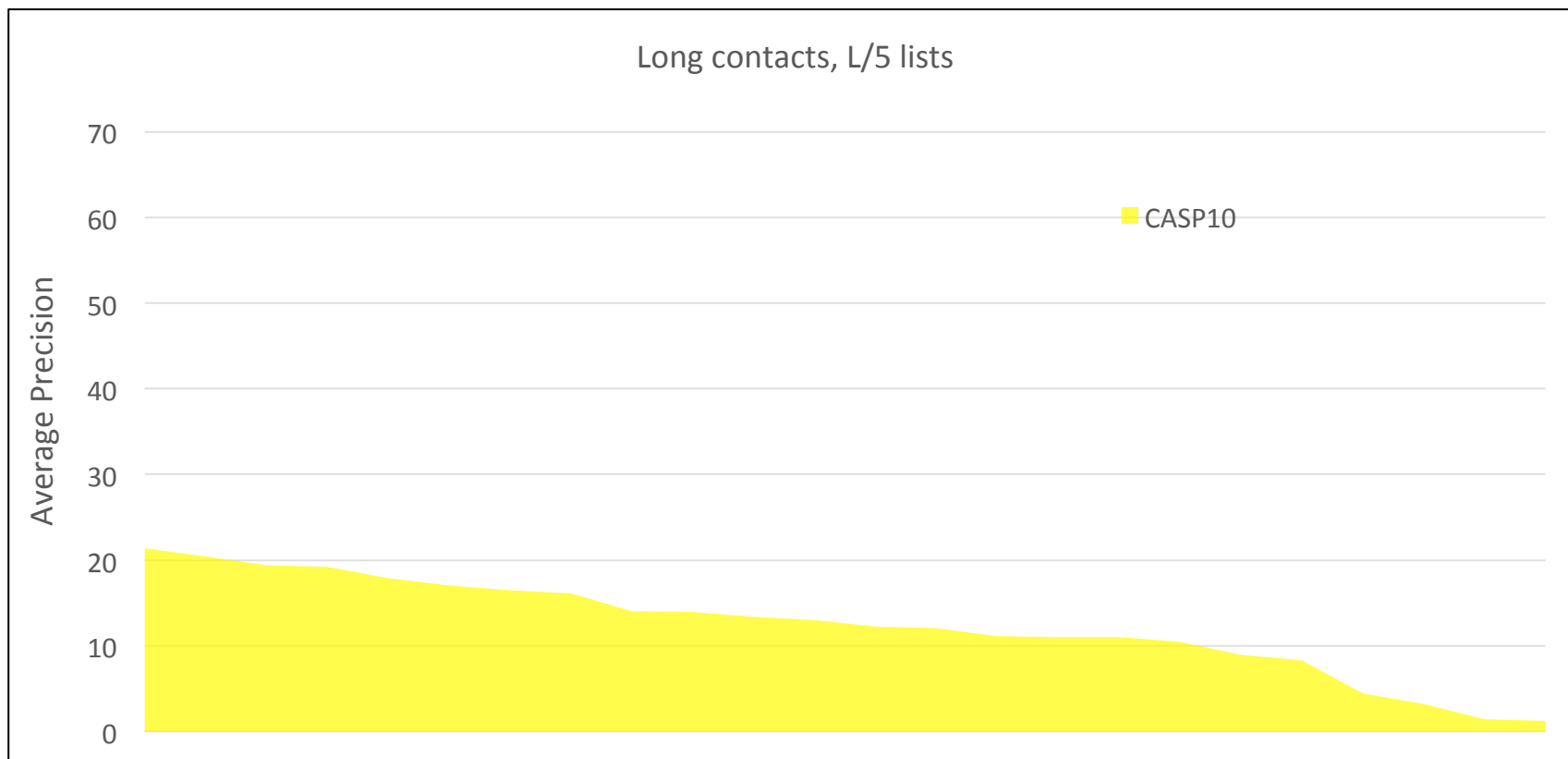
Long contacts, (FM only), sum Zscore (>0)



Extra long contacts only, (FM only), Zscore>0

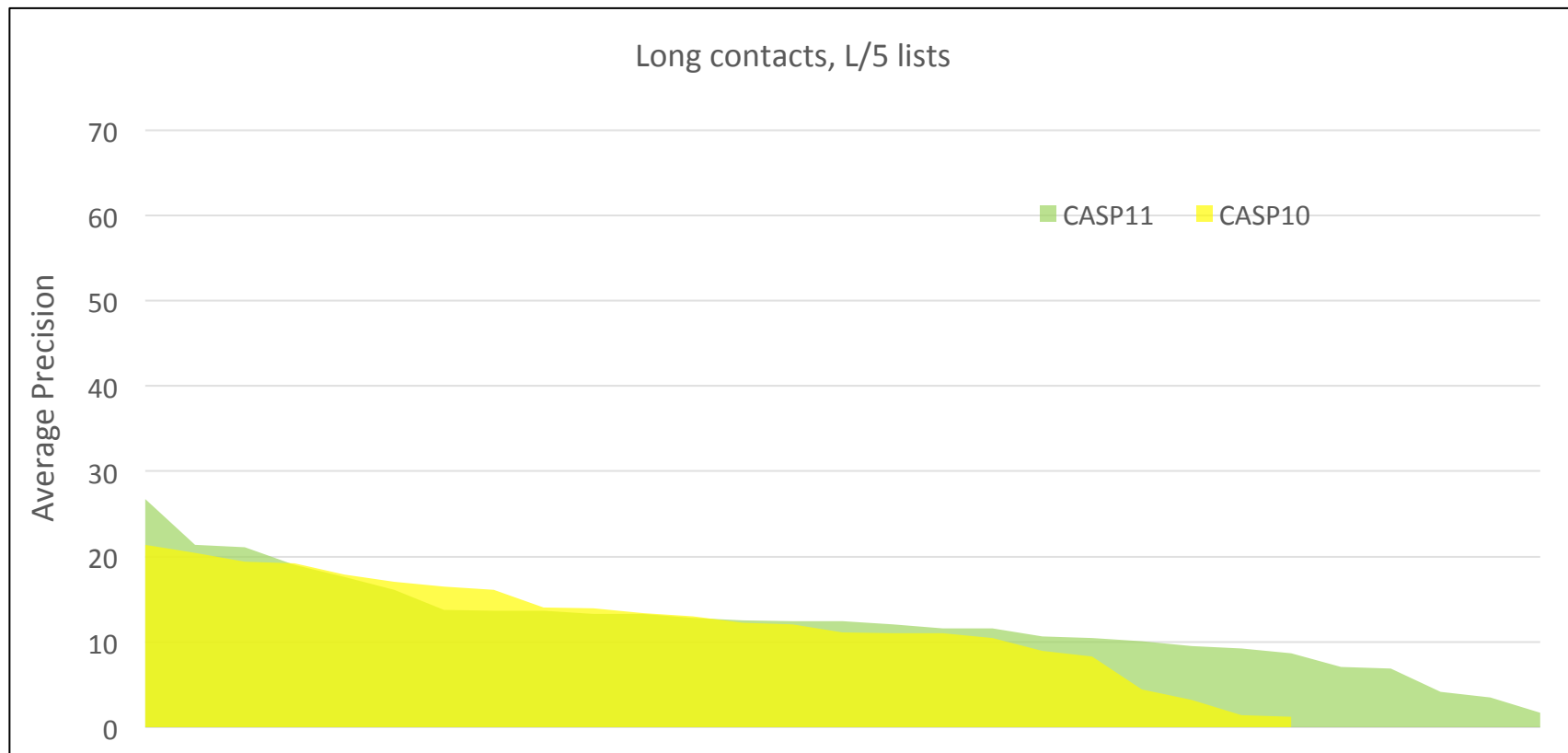


Improvement in contact prediction accuracy over CASP10-13 meetings



CASP10: 23 groups, 15 non-redundant

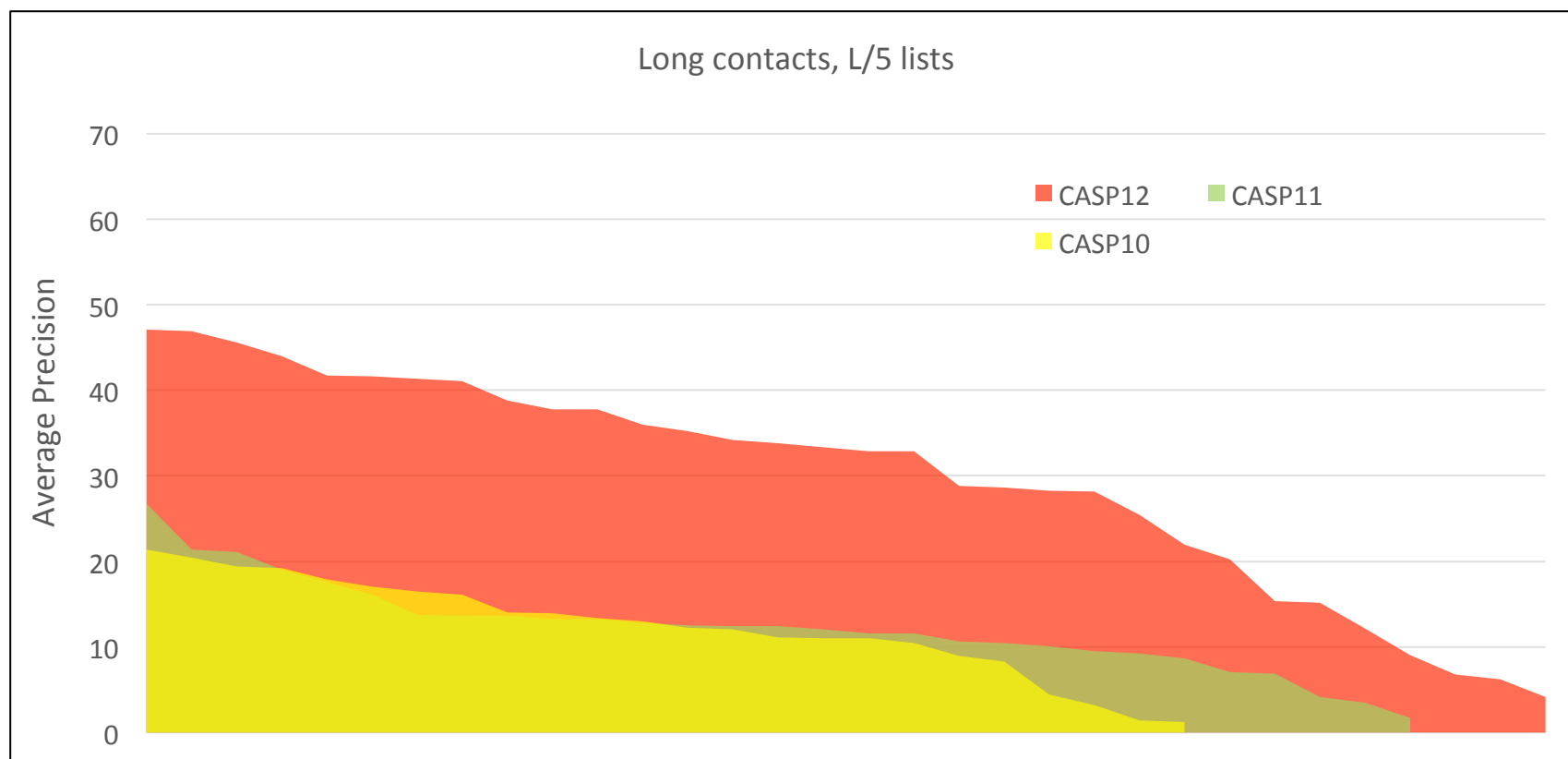
Improvement in contact prediction accuracy over CASP10-13 meetings



CASP10: 23 groups, 15 non-redundant

CASP11: 28 groups, 22 non-redundant

Improvement in contact prediction accuracy over CASP10-13 meetings

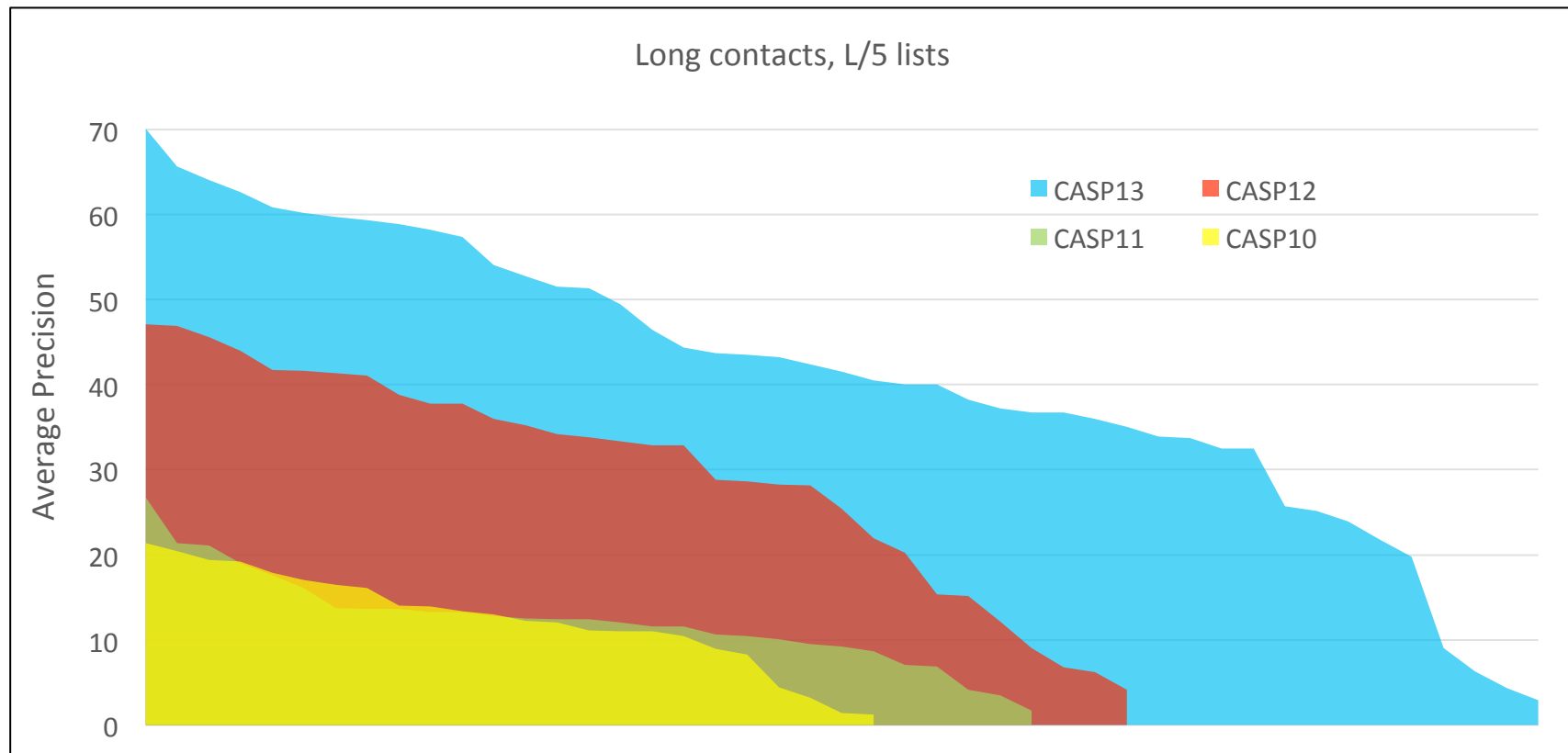


CASP10: 23 groups, 15 non-redundant

CASP11: 28 groups, 22 non-redundant

CASP12: 31 groups, 22 non-redundant

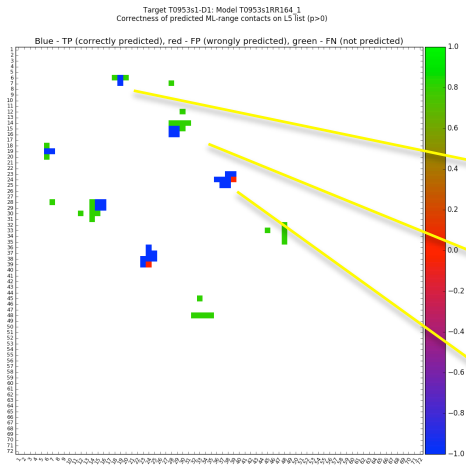
Improvement in contact prediction accuracy over CASP10-13 meetings



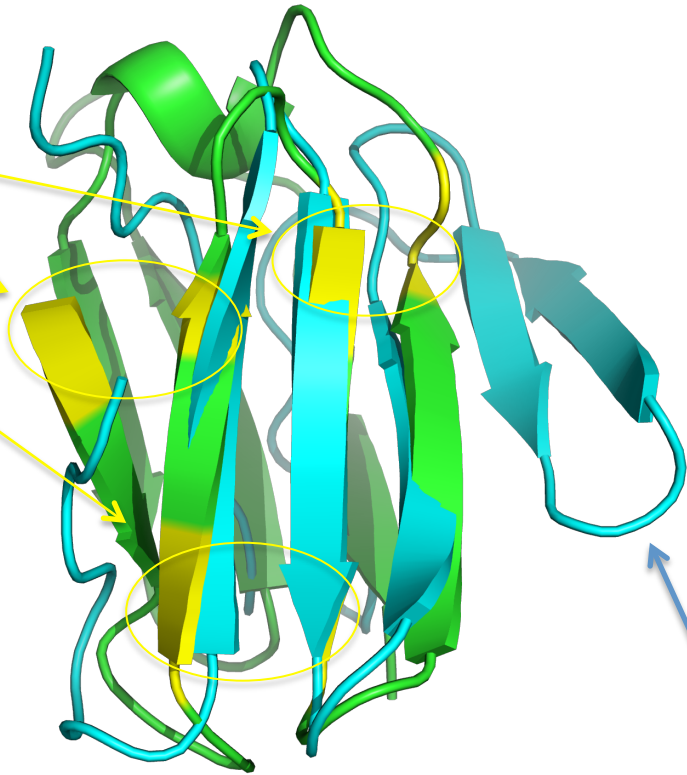
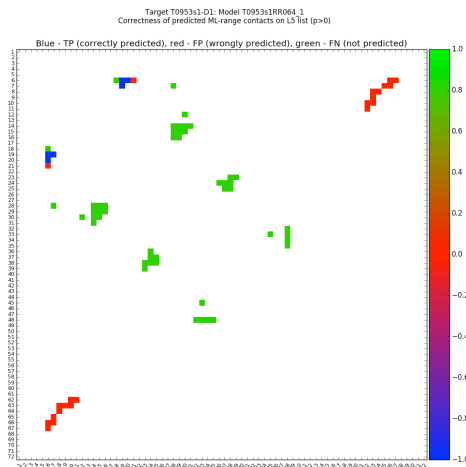
CASP10: 23 groups, 15 non-redundant
CASP11: 28 groups, 22 non-redundant
CASP12: 31 groups, 24 non-redundant
CASP13: 44 groups, 34 non-redundant

T0953s1d1

Good Fscore 63.4

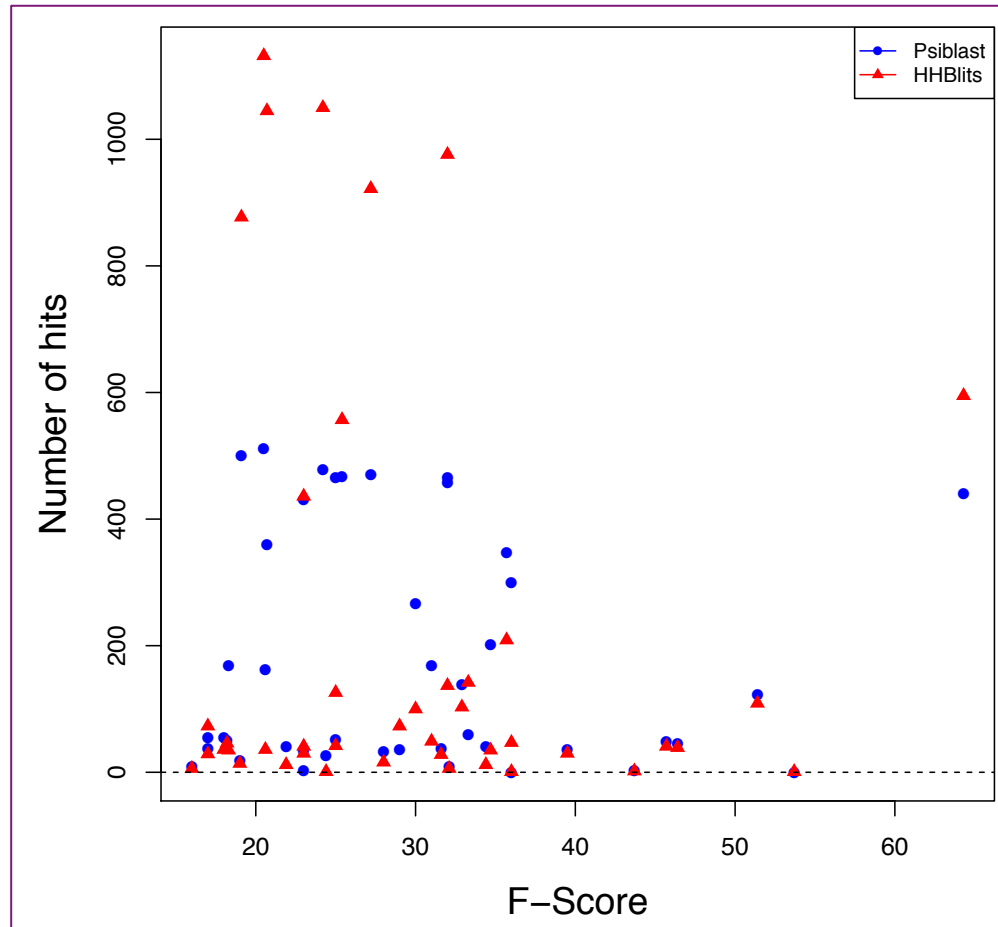


Poor Fscore: 14.64



Best TS model (G43), Cyan, GDT_TS 54.48
Contact model (G164), Green, GDT_TS 41.05

Relationship between *sequence profile depth* and success (F-score) of predicting contacts

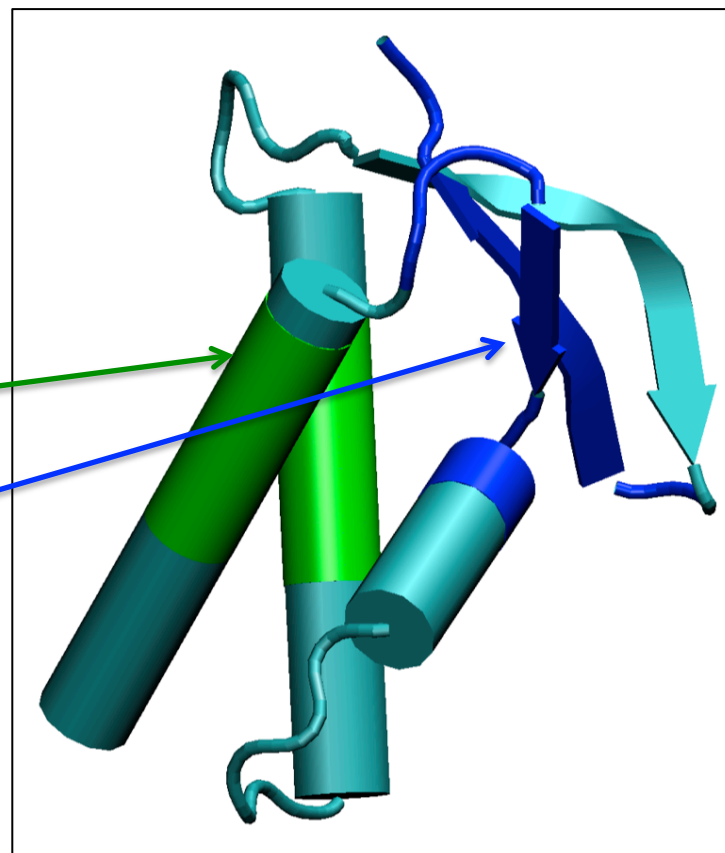
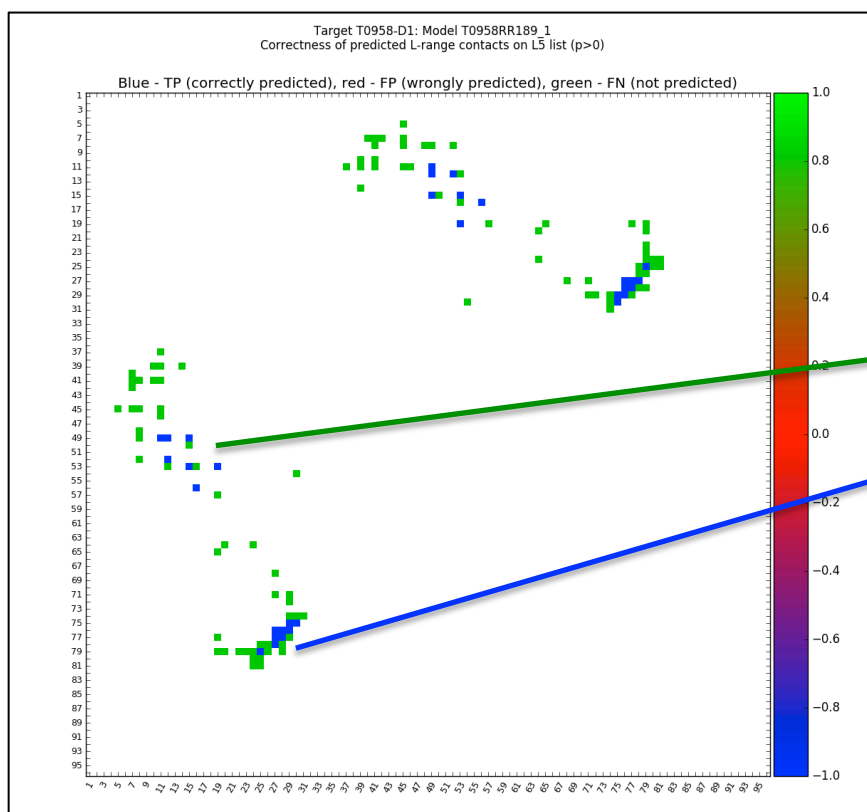


- Less reliant on sequence profiles.

Limited signal coming from sequence

	Blast			Blast+HHblits			
Fscore	e-5	e-20	Neff	Fscore	e-5	e-20	Neff
20.7	23	10	360	18	58	31	54
32.9	252	25	138	25	58	31	51
45.7	14	1	49	→ 53.7	1	1	0
17	46		37	29	14	14	36
28	50	17	33	→ 24.2	1266	1028	478
→ 31.6	40	22	37	→ 43.7	1	1	2
39.5	46	37	36	→ 19.1	3752		500
21.9	669	37	40	→ 32.1	4	3	9
34.4	591		41	23	584	110	430
31	89	46	168	→ 36	1	1	0
33.3	38	20	60	→ 19	7	4	18
20.5	3021	1905	511	→ 51.4	21	6	123
32	172	172	457	→ 51.4	77	13	123
25.4	609	183	467	30	231	126	267
25	6130	129	465	64.3	302	85	441
34.7	132	31	201	64.3	545	343	441
34.7	91	111	201	27.2	1730	68	470
16	30	6	9	18.2	629	1163	50
17	30	6	55	18.2	3755	38	50
35.7	278	17	347	32	1380	53	465
23	19	9	35				
20.6	38	23	162				
18.3	38	23	169				
36	194	1	300				
24.4	194	1	27				
46.4	58	31	45				

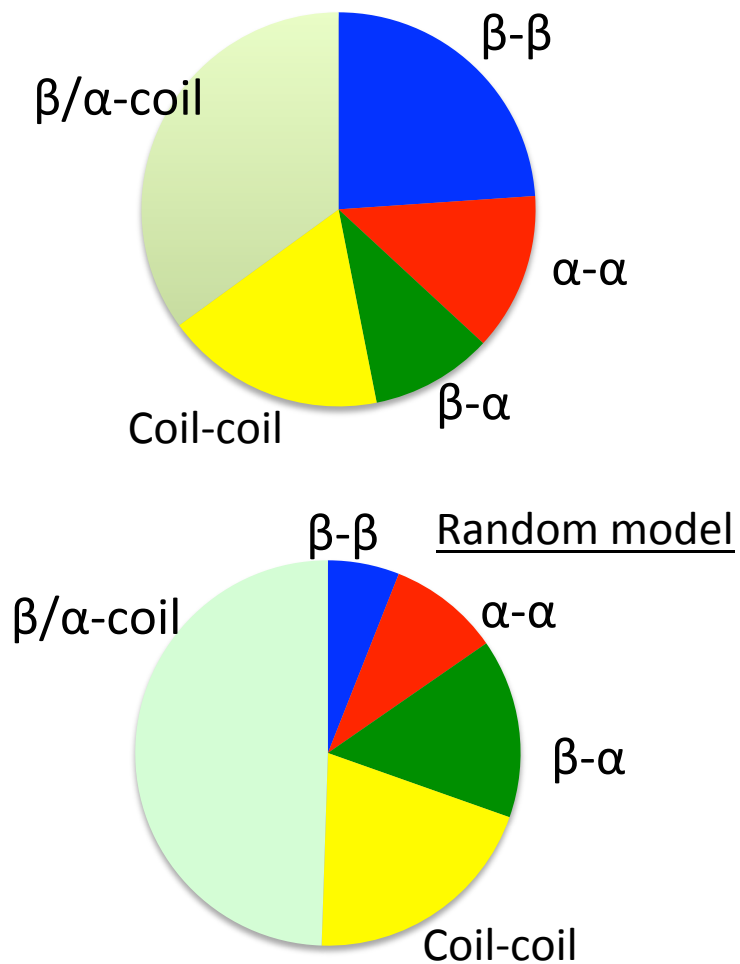
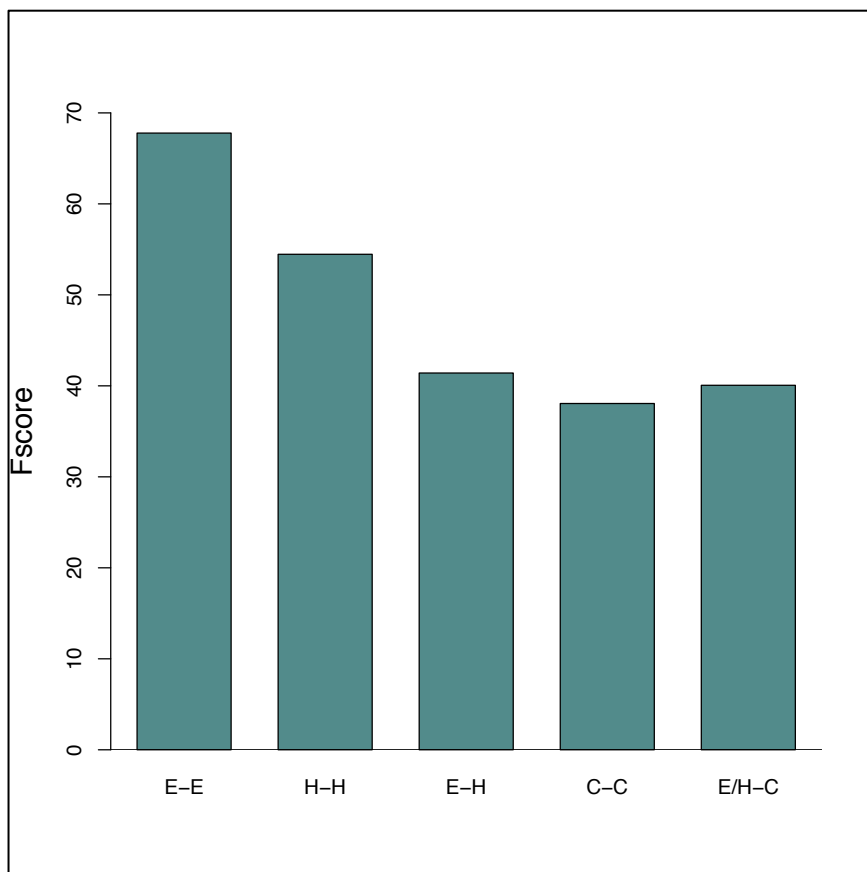
What is what?



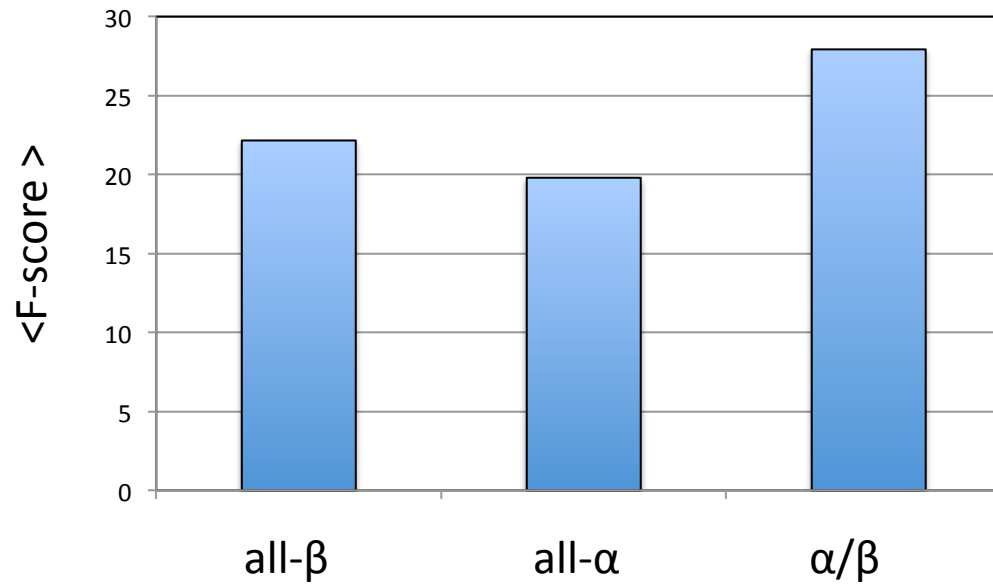
Green: parallel (parallel with diagonal) + diffuse (helical)

Blue: Anti-parallel (orthogonal with diagonal) + compact (strand)

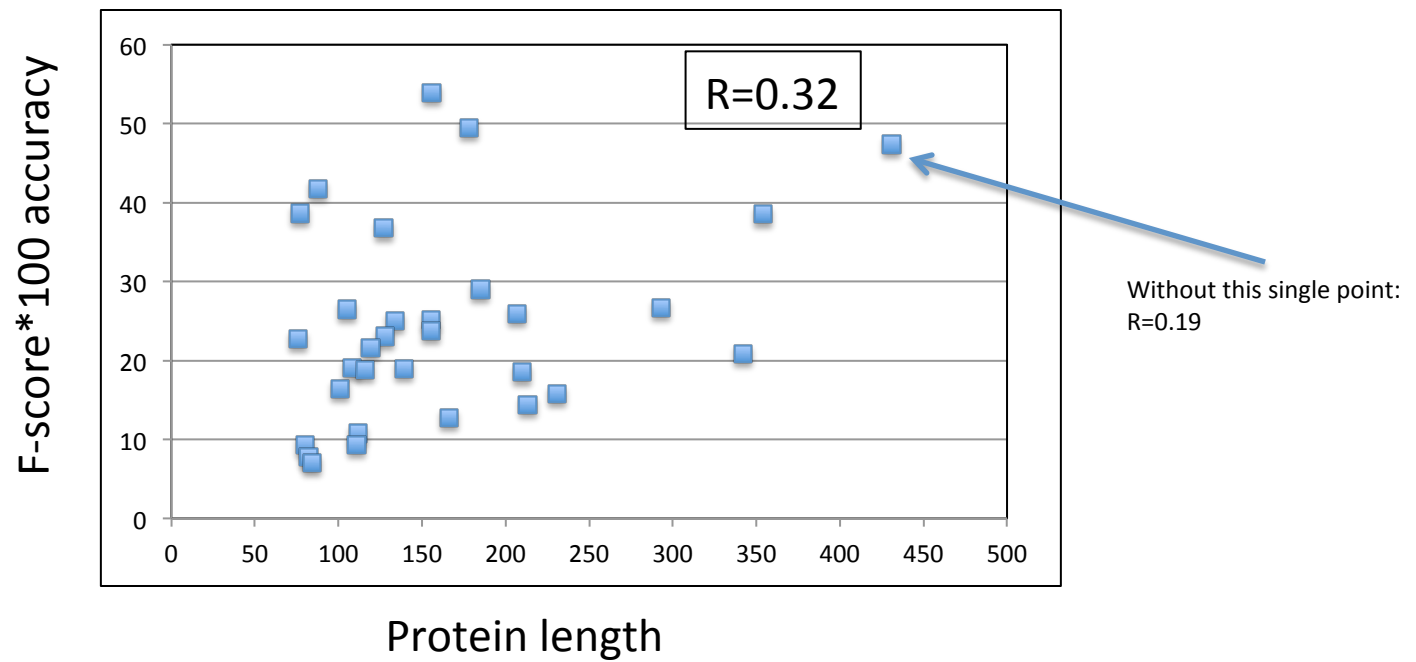
Performance vs. secondary structure interactions



Topology dependence of success rates, Class level



Correlation with size



Conclusions

- Contact prediction methods made a major advance for the last two years
- A lot of different subsets of correct contacts can be made and used successfully in 3D modeling
- It is difficult to directly correlate predicted contacts with 3D predictions because of ambiguity and lack of overlap between categories but :
 - Best 3D predictors have either even superior contact predictions or better ways to use contact information
 - From the few examples when both contacts and 3D structures were predicted we see strong inconsistencies: it is important to know how to use contact information
- Often very few homologous sequences were available, but very good contact predictions were made
 - Less emphasis on co-variance based methods (supported by the abstract of invited groups)

Acknowledgement

CASP and Predictioncenter at UC Davies, Davies, USA:

Andriy Kryshatafovych
Bohdan Monastyrskyy
Krzysztof Fidelis

CASP organizers

Albert Einstein College of Medicine, New York, USA:

Rojan Shrestha
Eduardo Fajardo
Nelson Gil