

## Supplementary Materials for **Protein Complex–Based Analysis Framework for High-Throughput Data Sets**

Arunachalam Vinayagam,\* Yanhui Hu, Meghana Kulkarni, Charles Roesel, Richelle Sopko, Stephanie E. Mohr, Norbert Perrimon\*

\*To whom correspondence should be addressed. E-mail: [vinu@genetics.med.harvard.edu](mailto:vinu@genetics.med.harvard.edu) (A.V.); [perrimon@receptor.med.harvard.edu](mailto:perrimon@receptor.med.harvard.edu) (N.P.)

Published 26 February 2013, *Sci. Signal.* **6**, rs5 (2013)  
DOI: 10.1126/scisignal.2003629

### **This PDF file includes:**

- Fig. S1. Schematic representation of protein complex scoring.
- Fig. S2. Snapshots of the COMPLEAT Web interface.
- Fig. S3. Complex enrichment results of baseline and EGF stimulus.
- Fig. S4. Baseline compared with EGF stimulus common complexes.
- Fig. S5. Baseline compared with EGF stimulus dynamic complexes: opposing effects.
- Fig. S6. Baseline compared with EGF stimulus: baseline-specific dynamic complexes.
- Fig. S7. Baseline compared with EGF stimulus: stimulus-specific dynamic complexes.
- Fig. S8. Complex enrichment results of baseline and insulin stimulus.
- Fig. S9. Baseline compared with insulin stimulus: common complexes.
- Fig. S10. Baseline compared with insulin stimulus: baseline-specific dynamic complexes.
- Fig. S11. Baseline compared with insulin stimulus: stimulus-specific dynamic complexes.
- Table S1. Compilation of literature protein complexes for humans, *Drosophila*, and yeast.
- Table S2. PPI data sets used to construct integrated PPI networks for humans, *Drosophila*, and yeast.
- Table S3. Predicted protein complexes for humans, *Drosophila*, and yeast.
- Table S4. Redundancy in the protein complex resource.
- Table S5. Overlap of the literature and predicted complexes at the protein level.
- Table S6. Proteome covered by the protein complex resources.

Table S7. Comparison of protein complexes with GO and KEGG with respect to co-citation.

Table S8. Comparison of protein complexes with GO and KEGG with respect to protein colocalization.

Table S9. Comparison of protein complexes with GO and KEGG with respect to gene coexpression.

Table S10. Annotation of the protein complex resource.

Table S11. Gene or protein input identifiers supported by the COMPLEAT.

Table S20. Dynamic phosphosites changing in response to insulin treatment.

**Other Supplementary Material for this manuscript includes the following:**

(available at [www.sciencesignaling.org/cgi/content/full/6/264/rs5/DC1](http://www.sciencesignaling.org/cgi/content/full/6/264/rs5/DC1))

Table S12 (Microsoft Excel format). Enriched protein complexes at baseline (mtDER-S2R+ cell line).

Table S13 (Microsoft Excel format). Enriched protein complexes at EGF stimulus (mtDER-S2R+ cell line).

Table S14 (Microsoft Excel format). Enriched protein complexes at baseline (S2R+ cell line).

Table S15 (Microsoft Excel format). Enriched protein complexes at insulin stimulus (S2R+ cell line).

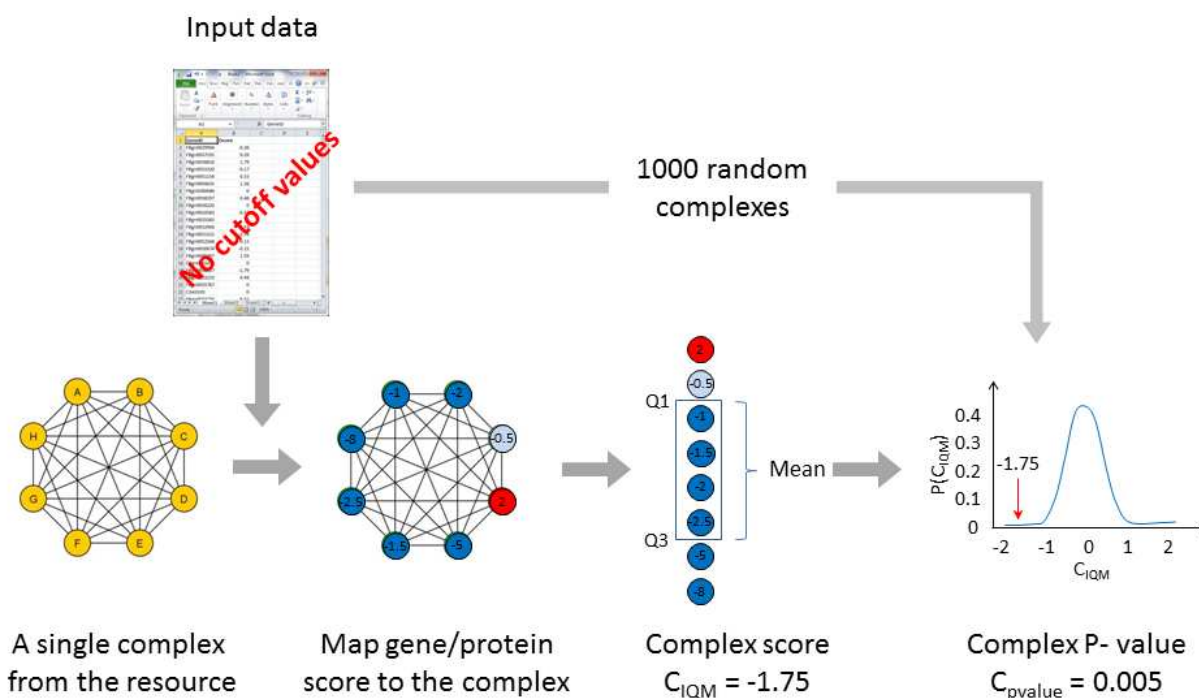
Table S16 (Microsoft Excel format). Consistent protein complexes with respect to baseline versus EGF stimulus.

Table S17 (Microsoft Excel format). Dynamic protein complexes with respect to baseline versus EGF stimulus.

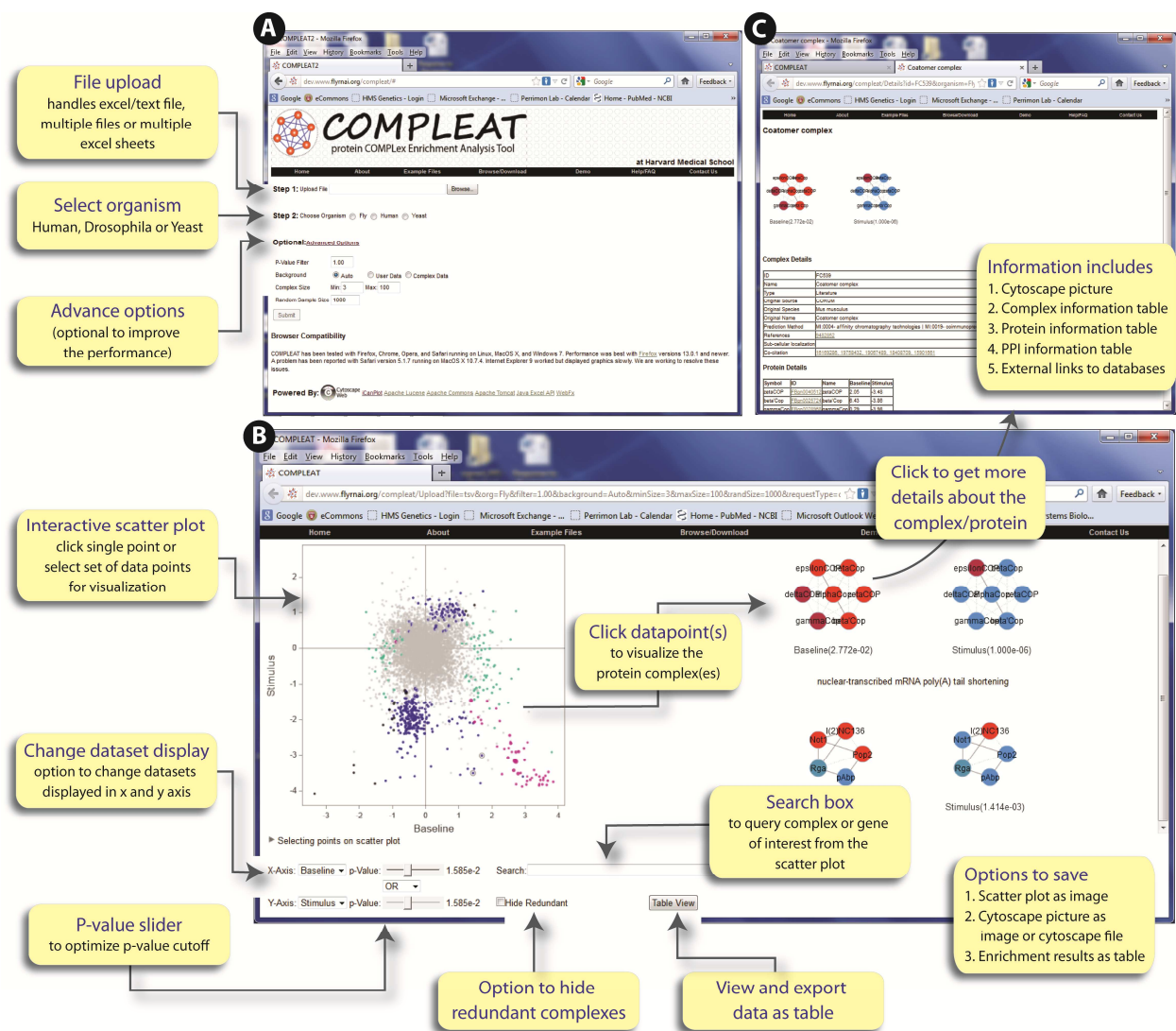
Table S18 (Microsoft Excel format). Consistent protein complexes with respect to baseline versus insulin stimulus.

Table S19 (Microsoft Excel format). Dynamic protein complexes with respect to baseline versus insulin stimulus.

# Supplementary materials

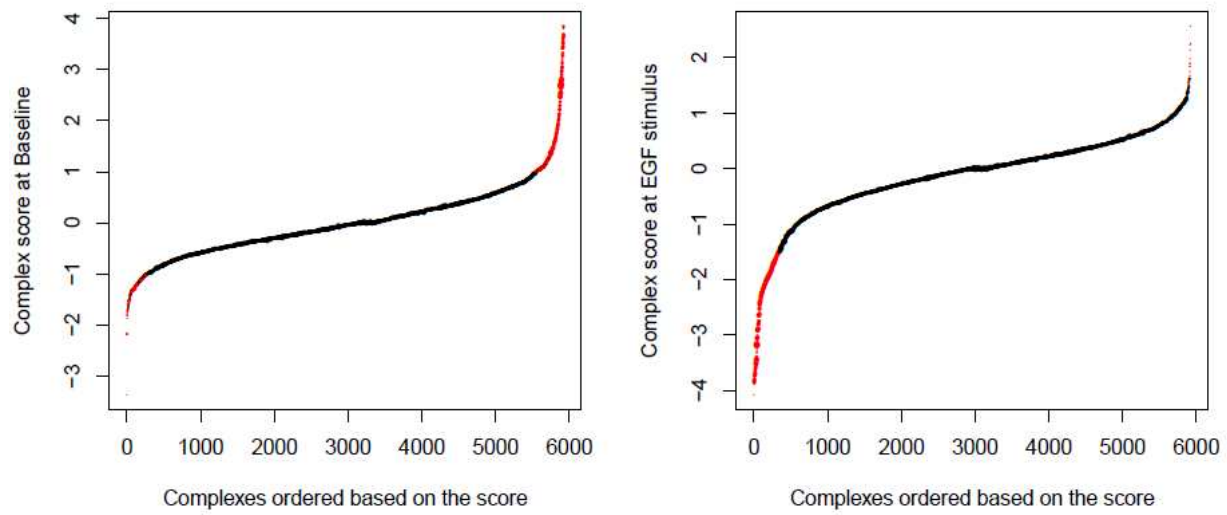


**Figure S1:** Schematic representation of protein complex scoring. In this model, the complex score and p-value calculation of a single complex is shown. First, the input data (without preselecting hits) is mapped to the protein complex. To calculate the interquartile mean (IQM), complex members are ordered based on the protein-score, and the mean value between first (Q1) and third (Q3) quartile is calculated. The p-value corresponding to the IQM is calculated by comparing it to the distribution of random IQM scores calculated based on the 1000 random complexes. Random complexes are generated either based on the input data or based on the complex resource, depending on the user specification.

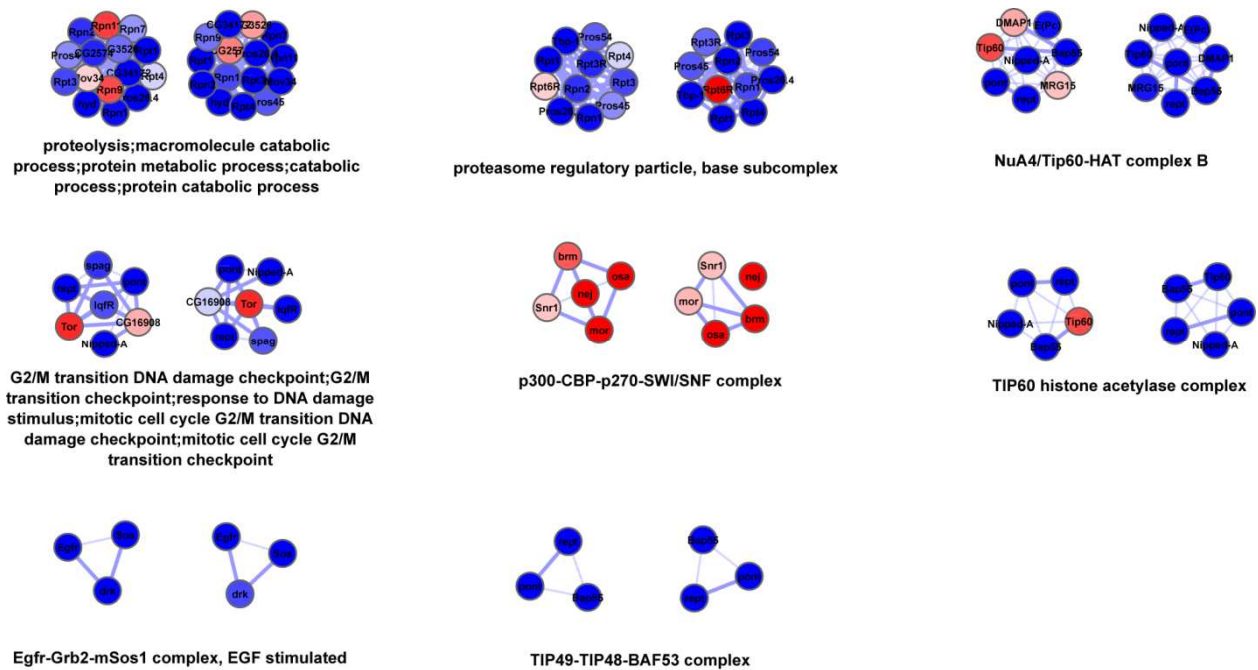


**Figure S2:** Snapshots of the COMPLEAT Web interface. (A) Input page for COMPLEAT with options to upload input file, choose organism and set advance parameters. (B) The COMPLEAT result page includes an interactive scatterplot where each point on the scatter plot represents a single complex whose position corresponds to the score. Size reflects the relative complex size, and color corresponds to the p-value. The user has the option to change the p-value threshold using p-value adjustment sliders. When a user selects the complex of interest from the scatter-plot, the network illustrations of the complexes are displayed on the Web Cytoscape panel (right panel of the same page). The node color in the network corresponds to the user input values, and the color-code ranges from blue to red (blue corresponds to the lowest value, and red is the maximum value). Note that the gray node represents a missing value, meaning that a particular gene or protein is present in the complex but missing in the user input data.

There are two types of edges: Solid edges correspond to known PPIs. Broken edges are interologs (proteins for which the ortholog gene pairs in another species are known to physically interact). The user has the options to zoom in or out in the network and save the network images. (C) Additional information about complexes or proteins can be obtained by clicking nodes or complexes. For example, clicking a node takes the user to the corresponding gene or protein database. Clicking a complex provides annotation regarding the complex, such as the original source, purification method or prediction algorithm, PubMed references (if available), sub-cellular locations and co-cited literature (see Materials and Methods for details).



**Figure S3:** Complex enrichment results of baseline and EGF stimulus. (A) Distribution of complex scores from baseline data. Significant complexes are highlighted in red ( $p\text{-value} \leq 0.01$  and score  $\geq 1$  or  $\leq -1$ ). (B) Complex score distribution from EGF stimulus data. Significant complexes are shown in red ( $p\text{-value} \leq 0.01$  and the score  $\geq 1.5$  or  $\leq -1.5$ ). The point size is proportional to the complex size.



**Figure S4:** Baseline compared with EGF stimulus common complexes. Non-redundant complexes corresponding to table S16 are shown. Each complex is represented twice; the complex on the left corresponds to baseline, and that to the right represents the stimulus condition. The network picture was generated using Cytoscape software ([www.cytoscape.org/](http://www.cytoscape.org/)). The node color ranges from dark blue to dark red, where the lowest value correspond to dark blue (negative Z-score) and highest score corresponds to dark red (positive z-scores). Solid edges correspond to known PPI and broken edges correspond to interolog.

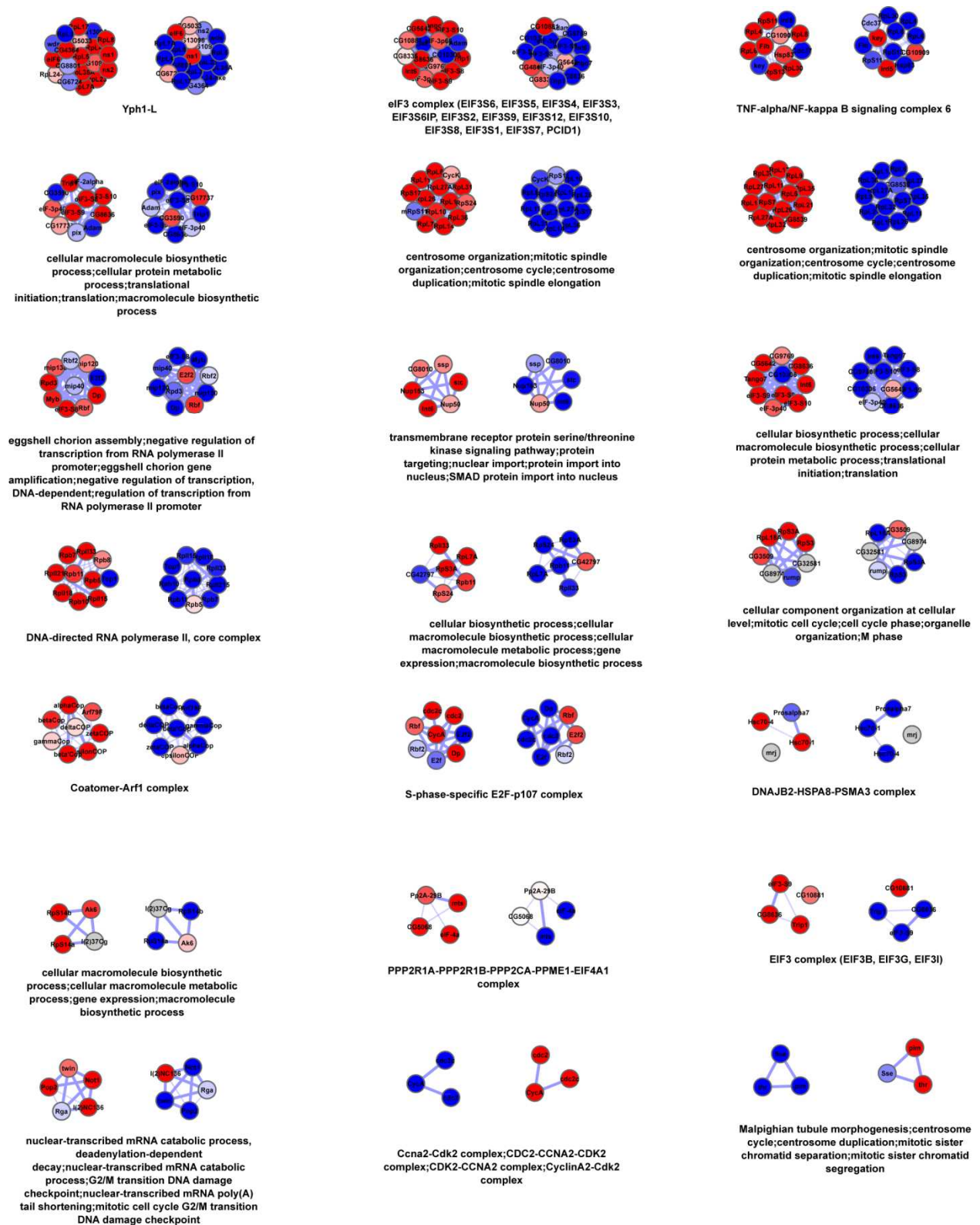
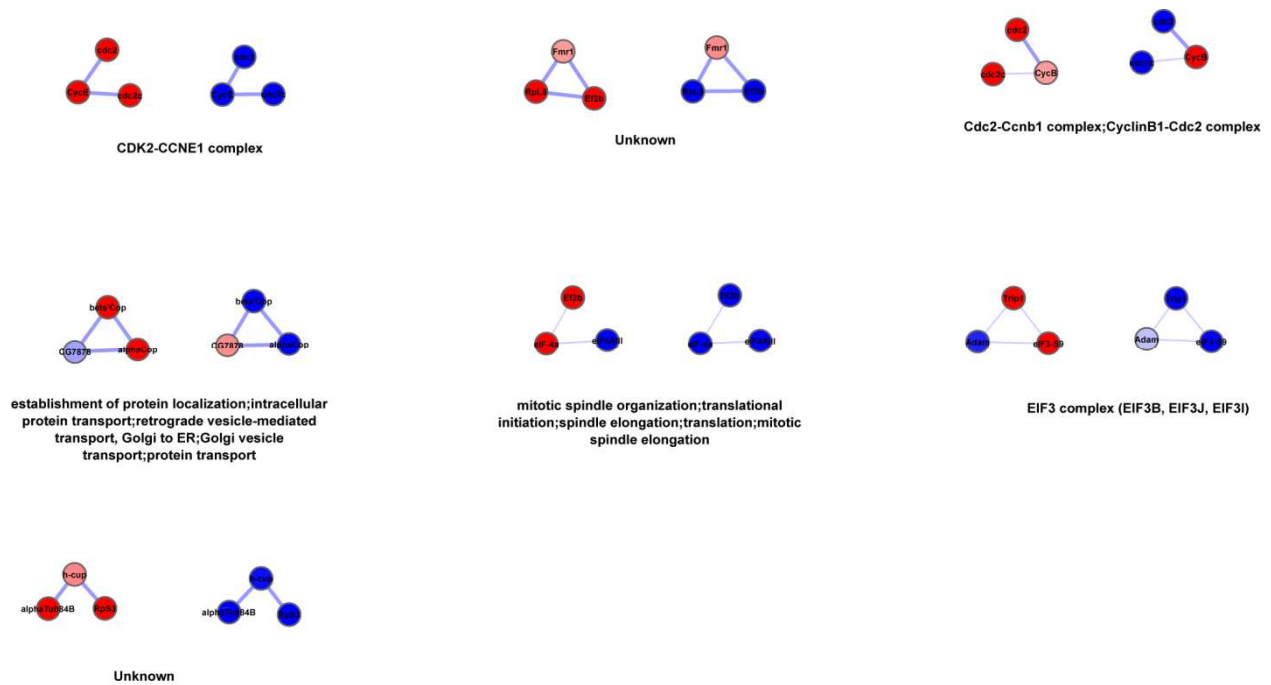


Figure S5: continued...



**Figure S5:** Baseline compared with EGF stimulus dynamic complexes: opposing effects. Non-redundant complexes corresponding to table S17 are shown. Node color and edge style are as described in figure S4.



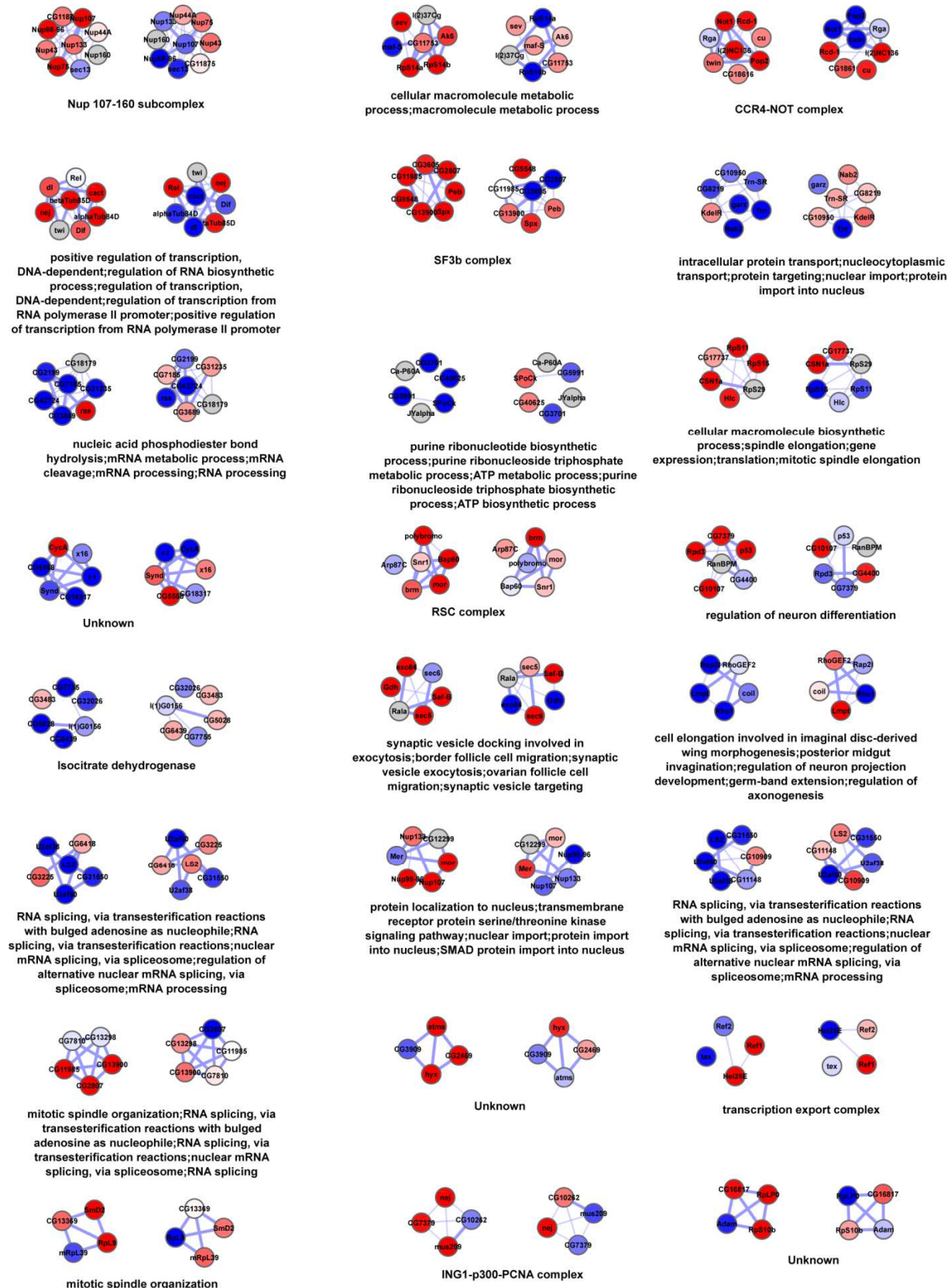
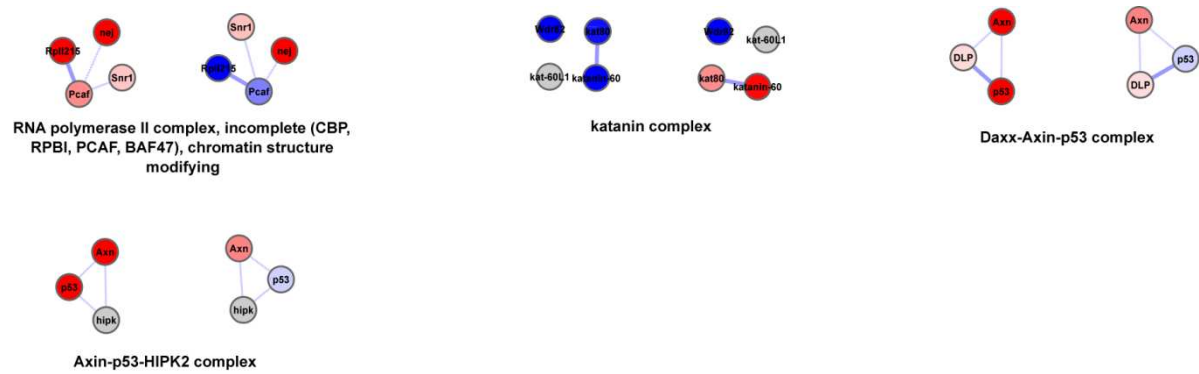


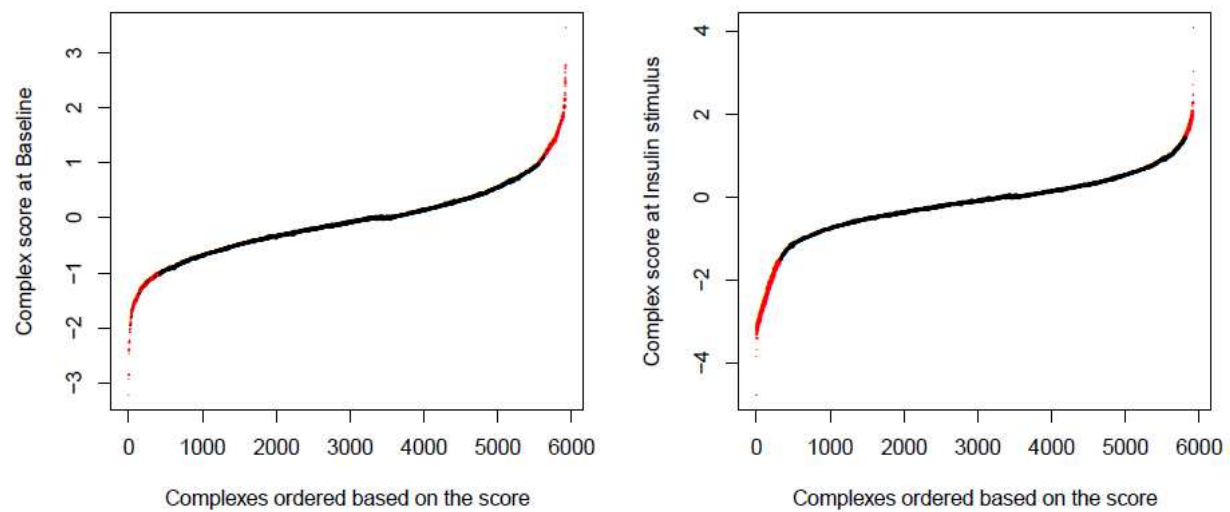
Figure S6: Continued...



**Figure S6:** Baseline compared with EGF stimulus: baseline-specific dynamic complexes. Non-redundant complexes corresponding to table S17 are shown. Node color and edge style are as described in figure S4.







**Figure S8:** Complex enrichment results of baseline and insulin stimulus. (A) Baseline complex scores distribution. Significant complexes are highlighted in red ( $p\text{-value} \leq 0.01$  and score  $\geq 1$  or  $\leq -1$ ). (B) Complex score distribution from insulin stimulus data. Significant complexes are shown in red ( $p\text{-value} \leq 0.01$  and score  $\geq 1.5$  or  $\leq -1.5$ ). The point size is proportional to the complex size.

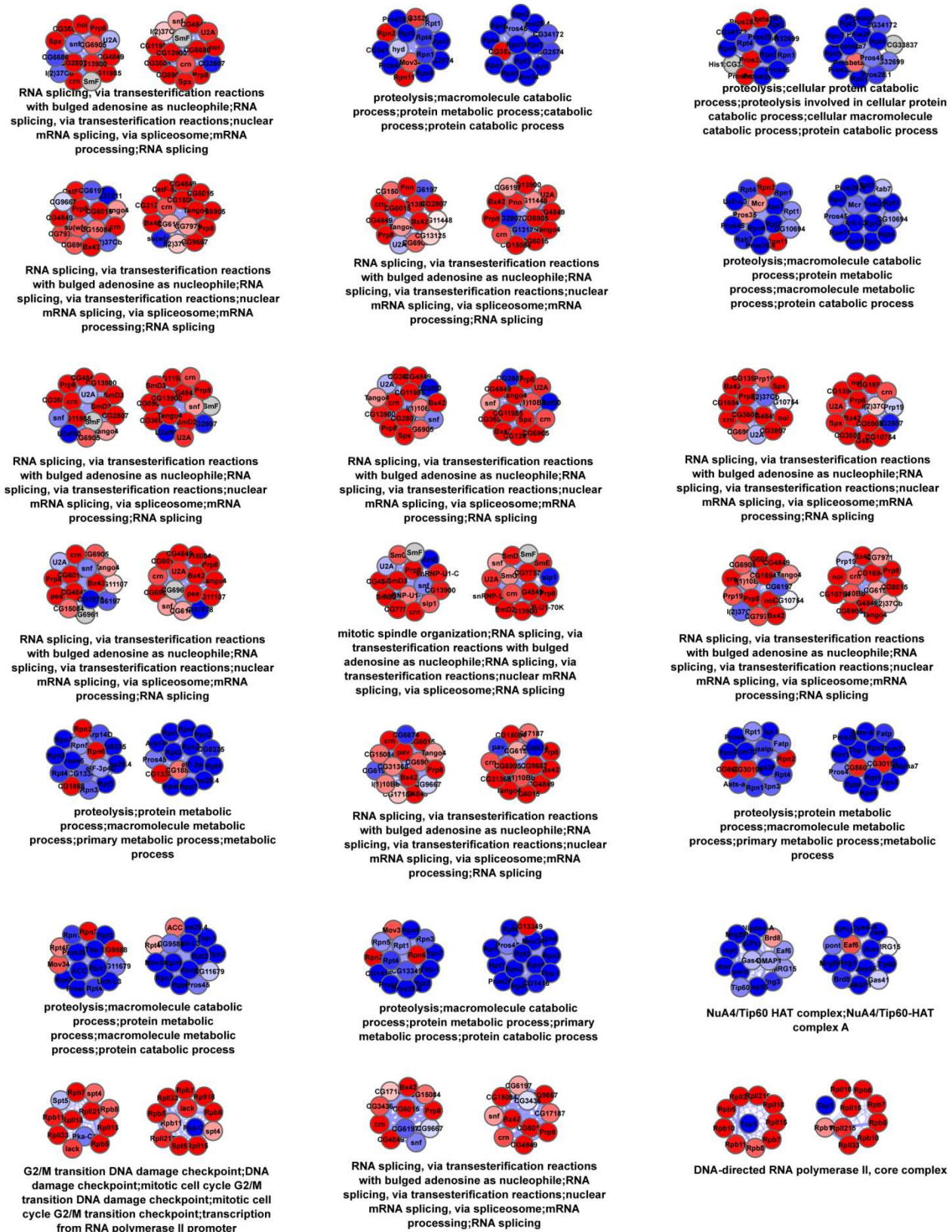
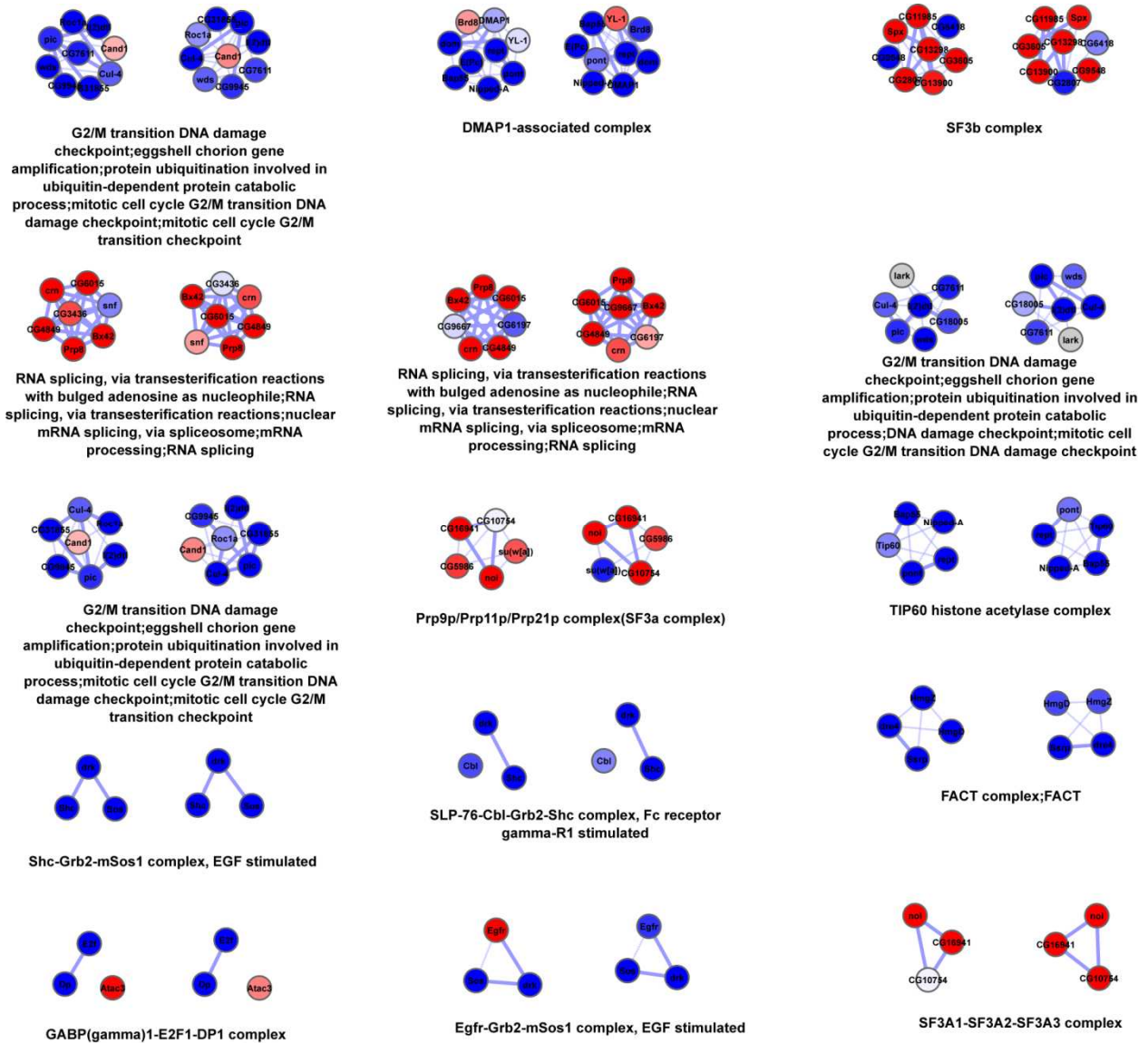


Figure S9: Continued...



**Figure S9:** Baseline compared with insulin stimulus: common complexes. Non-redundant complexes corresponding to table S18 are shown. Node color and edge style are as described in figure S4.

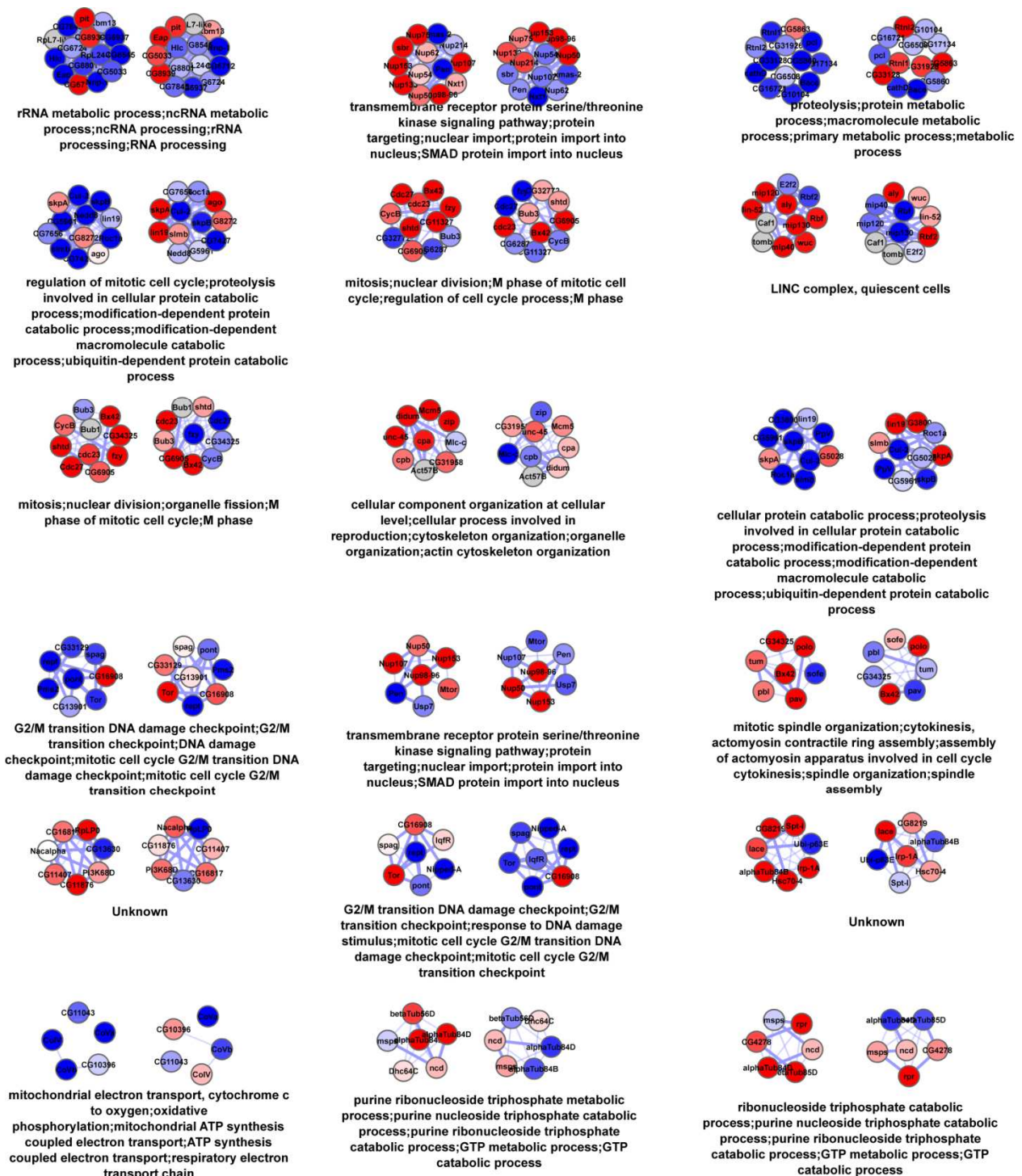
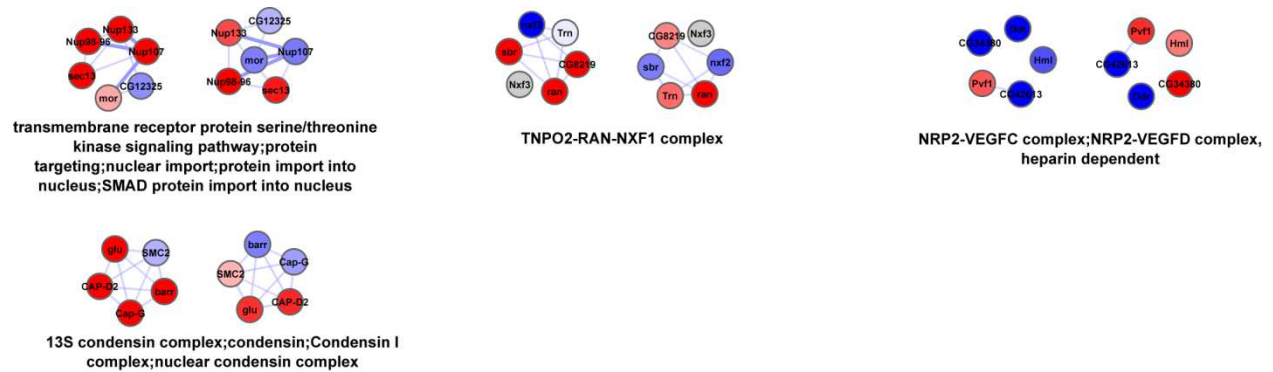
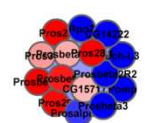


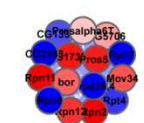
Figure S10: Continued...



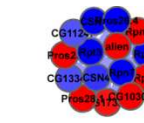
**Figure S10:** Baseline compared with insulin stimulus: baseline-specific dynamic complexes. Non-redundant complexes corresponding to table S19 are shown. Node color and edge style are as described in figure S4.



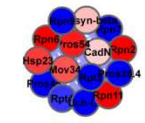
cellular protein catabolic process;proteolysis involved in cellular protein catabolic process;modification-dependent protein catabolic process;modification-dependent macromolecule catabolic process;ubiquitin-dependent protein catabolic process



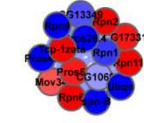
proteolysis;protein metabolic process;macromolecule metabolic process;primary metabolic process;metabolic process



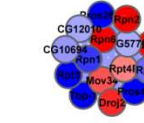
proteolysis;cell cycle phase;protein metabolic process;response to DNA damage stimulus;protein catabolic process



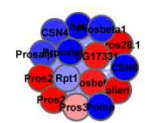
proteolysis;protein metabolic process;macromolecule catabolic process;primary metabolic process;protein catabolic process



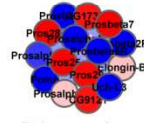
cellular response to stress;proteolysis;protein metabolic process;cellular macromolecule metabolic process;response to DNA damage stimulus



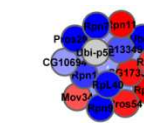
cellular response to stress;proteolysis;macromolecule catabolic process;protein catabolic process;response to DNA damage stimulus



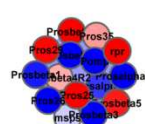
proteolysis;cellular protein catabolic process;proteolysis involved in cellular protein catabolic process;response to DNA damage stimulus;protein catabolic process



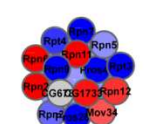
cellular protein catabolic process;proteolysis involved in cellular protein catabolic process;modification-dependent protein catabolic process;modification-dependent macromolecule catabolic process;ubiquitin-dependent protein catabolic process



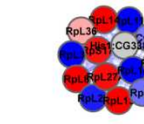
proteolysis;cellular protein catabolic process;proteolysis involved in cellular protein catabolic process;cellular macromolecule catabolic process;protein catabolic process



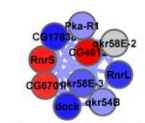
cellular protein catabolic process;proteolysis involved in cellular protein catabolic process;modification-dependent protein catabolic process;modification-dependent macromolecule catabolic process;ubiquitin-dependent protein catabolic process



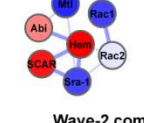
proteolysis;protein metabolic process;macromolecule metabolic process;primary metabolic process;metabolic process



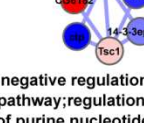
mitotic spindle organization;centrosome organization;centrosome cycle;centrosome duplication;mitotic spindle elongation



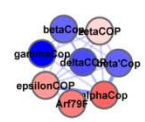
activation of cysteine-type endopeptidase activity involved in apoptotic process;regulation of cysteine-type endopeptidase activity involved in apoptotic process;activation of cysteine-type endopeptidase activity;positive regulation of cysteine-type endopeptidase activity;positive regulation of cysteine-type endopeptidase activity involved in apoptotic process



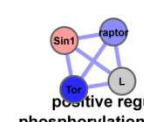
Wave-2 complex (Rac-activated)



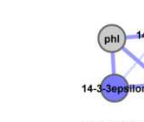
negative regulation of insulin receptor signaling pathway;regulation of GTPase activity;regulation of purine nucleotide catabolic process;regulation of GTP catabolic process;regulation of nucleotide catabolic process



Coatamer-Arf1 complex

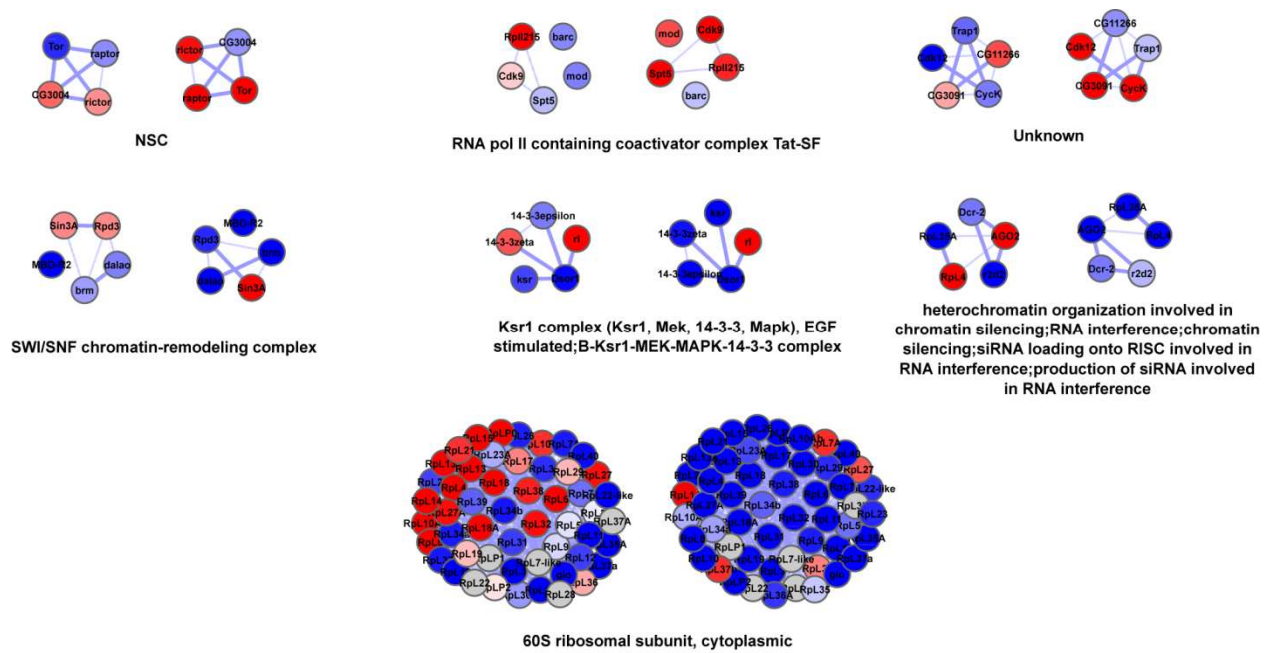


positive regulation of protein phosphorylation;regulation of protein phosphorylation;positive regulation of protein modification process;positive regulation of phosphorylation;regulation of protein modification process



BRAF-MAP2K1-MAP2K2-YWHA complex;RAF1-MAP2K1-YWHA complex

Figure S11: Continued...



**Figure S11:** Baseline compared with insulin stimulus: stimulus-specific dynamic complexes. Non-redundant complexes corresponding to table S19 are shown. Node color and edge style are as described in figure S4.

**Table S1:** Compilation of literature protein complexes for humans, *Drosophila*, and yeast. Complex source are either from literature curation (LC) or high confidence complexes reported in literature based on high-throughput MS-pull down data (HT). Source organism: the organism in which the protein complex is reported in the database or publication. Ortholog mapping: the protein complexes were mapped to other organism using DIOPT (17), an ortholog mapping tool.

Database/dataset	Source	Source organism	Ortholog mapping	Human	Drosophila	Yeast
CORUM (15)	LC	Human, mouse	Human, <i>Drosophila</i> , yeast	2363	2162	1395
PINdb (13)	LC	Human, yeast	Human, <i>Drosophila</i> and yeast	286	280	276
CYC2008 (14)	LC	Yeast	Human, <i>Drosophila</i>	358	346	408
	HT	Yeast	Human, <i>Drosophila</i>	343	333	400
Gene Ontology (2)	LC	Human, <i>Drosophila</i> , yeast	No mapping	282	146	304
<i>Drosophila</i> AP-MS pull-down complexes (16)	HT	<i>Drosophila</i>	Human, yeast	511	556	331
KEGG module (3)	LC	Human	<i>Drosophila</i> , yeast	210	196	193
Signalink (32)	LC	Human, <i>Drosophila</i>	yeast	14	14	14
FlyReactome (52)	LC	<i>Drosophila</i>	Human, yeast	9	9	9
All literature complex				3638	3077	2173

**Table S2:** PPI data sets used to construct integrated PPI networks for humans, *Drosophila*, and yeast. Name of the data set, publication reference, URL, number of PPIs and proteins in the data set are given. All the PPI datasets are downloaded from the corresponding Website and the database version corresponds to the March 2012 release. For humans and yeast, the protein or gene identifiers are mapped to NCBI Entrez gene identifier. In case of *Drosophila*, the gene or protein identifiers are mapped to Flybase gene identifier.

Database/datasets	Human		Drosophila		Yeast	
	PPIs	Proteins	PPIs	Proteins	PPIs	Proteins
BioGrid (33) <a href="http://thebiogrid.org/">http://thebiogrid.org/</a>	59226	12529	23916	7305	60062	5374
IntAct (53) <a href="http://www.ebi.ac.uk/intact/">http://www.ebi.ac.uk/intact/</a>	40368	9468	25385	7530	76147	5555
DIP (54) <a href="http://dip.doe-mbi.ucla.edu/dip/Main.cgi">http://dip.doe-mbi.ucla.edu/dip/Main.cgi</a>	2860	1855	19753	6584	22028	4675
MINT (55) <a href="http://cbm.bio.uniroma2.it/mint/">http://cbm.bio.uniroma2.it/mint/</a>	19048	6315	17336	5917	27680	5104
HPRD (35) <a href="http://www.hprd.org/">http://www.hprd.org/</a>	39172	9670	-	-	-	-
DroID (36) <a href="http://www.droidb.org/">http://www.droidb.org/</a>	-	-	87070	9068	-	-
Drosophila AP-MS dataset (16)	-	-	10964	2296	-	-
MasterNet (integrated network)	108059	14495	98500	9373	118603	5729

**Table S3:** Predicted protein complexes for humans, *Drosophila*, and yeast. CFinder (38) was applied for humans, *Drosophila* and yeast filtered PPI networks. NetworkBLAST (19) was applied to identify protein complexes by aligning human and *Drosophila*, human and yeast, and *Drosophila* and yeast PPI networks.

Prediction source	Human		Drosophila		Yeast	
	Complexes	Proteins	Complexes	Proteins	Complexes	Proteins
CFinder / Human PPI	713	2046	-	-	-	-
CFinder / <i>Drosophila</i> PPI	-	-	433	1419	-	-
CFinder / Yeast PPI	-	-	-	-	423	1465
NetworkBLAST / Human vs. <i>Drosophila</i> PPI	1722	2665	1722	2638	-	-
NetworkBLAST / Human vs. Yeast PPI	3820	4369	-	-	3820	2712
NetworkBLAST / <i>Drosophila</i> vs. Yeast PPI	-	-	1532	2279	1532	1840
All predicted complexes	6251	6334	3639	3933	5551	3366

**Table S4:** Redundancy in the protein complex resource. A protein complex is defined as redundant if it is a subset, superset or shares 80% of proteins with the other complexes in the resource. Non-redundant complexes are constructed at 80%, meaning that no two complexes share more than 80% similarity. This table corresponds to Figure 2B.

Complexes	Human	Drosophila	Yeast
All complexes	9881	6703	7713
Non-redundant	5164	3399	3183
Redundant complexes	4717	3304	4530

**Table S5:** Overlap of the literature and predicted complexes at the protein level.

Organism	Literature specific		Predicted specific		Common		Total Proteins
	Proteins	Percentage	Proteins	Percentage	Proteins	Percentage	
<b>Human</b>	2959	31.8%	1769	19%	4565	49.1%	9293
<b>Drosophila</b>	2603	39.8%	917	14%	3016	46.1%	6536
<b>Yeast</b>	628	15.7%	714	17.9%	2652	66.4%	3994

**Table S6:** Proteome covered by the protein complex resources. This table corresponds to the Figure 2E.

Complexes	Human		Drosophila		Yeast	
	Covered	Total	Covered	Total	Covered	Total
All proteins	9293	20402	6536	13776	3994	5882
Conserved proteins	5161	6203	4009	5249	3184	3551

**Table S7:** Comparison of protein complexes with GO and KEGG with respect to co-citation. Significant and total protein complexes, GO categories and KEGG pathways are shown. Significant protein complexes, GO, and KEGG refers to significantly co-cited protein complexes compared to 1000 random sets of the same size ( $p < 0.05$ ). This table corresponds to Figure 2G.

Resource	Human		Drosophila		Yeast	
	Significant	Total	Significant	Total	Significant	Total
Complexes	8757	9125	4545	5817	6708	7098
GO	5119	5994	2170	2457	1901	2036
KEGG	100	102	110	129	167	178

**Table S8:** Comparison of protein complexes with GO and KEGG with respect to protein colocalization. Significant and total protein complexes, GO categories and KEGG pathways are shown. Significant protein complexes, GO, and KEGG refers to significantly colocalized protein complexes compared to 1000 random sets of the same size ( $p < 0.05$ ). This table corresponds to Figure 2H.

Resource	Yeast	
	Significant	Total
Complexes	3518	6682
GO	511	1976
KEGG	77	177

**Table S9:** Comparison of protein complexes with GO and KEGG with respect to gene coexpression. Significant and total protein complexes, GO categories and KEGG pathways are shown. Significant protein complexes, GO, and KEGG refers to significantly coexpressed genes compared to 1000 random sets of the same size ( $p < 0.05$ ). This table corresponds to Figure 2H.

Resource	Human		Drosophila	
	Significant	Total	Significant	Total
Complexes	5364	8450	4015	5491
GO	2082	6293	1787	2628
KEGG	126	175	159	189

**Table S10:** Annotation of the protein complex resource. Literature annotation corresponds to the annotation from the source database. GO enrichment was performed if the complex was predicted or if the annotation was available from the literature. Unknown complexes are new complexes for which no functional theme is associated. This table corresponds to Figure 2I.

Complexes	Human	Drosophila	Yeast
Literature annotation	3784	3372	2838
GO enrichment	5721	2614	4613
Unknown	376	717	262
Total	9881	6703	7713

**Table S11:** Gene or protein input identifiers supported by the COMPLEAT.

Identifier type	Human	Drosophila	Yeast
Symbol	Entrez gene symbol	Flybase gene symbol	Entrez gene symbol
Gene identifier	Entrez gene identifier	Entrez gene identifier	Entrez gene identifier
Protein identifier	Uniprot identifier	Uniprot identifier	Uniprot identifier
Species specific identifier	Entrez gene identifier	Flybase gene identifier	Locus tag

**Table S20:** Dynamic phosphosites changing in response to insulin treatment. A systematic investigation of insulin-induced phosphorylation using mass spectrometry and isobaric labeling of S2R+ cells identified dynamic phosphorylation of Moira and MBD-R2 following a 10-minute insulin stimulus including the Akt/RSK/S6 consensus motifs on Moira. This observation is consistent with human data, where Akt phosphorylates the human ortholog of Moira (BAP155) (48). Method: Biological duplicates of two conditions (no treatment or 10 minutes insulin treatment) were analyzed. Cells were lysed in 8 M urea, 75 mM NaCl, 50 mM Tris, pH 8.2, protease inhibitors cocktail (Roche), 1 mM NaF, 1 mM  $\beta$ -glycerophosphate, 1 mM sodium orthovanadate, 10 mM sodium pyrophosphate, 1 mM PMSF. 1mg of protein from each replicate was digested with trypsin (Promega) and processed as reported by Dephore and Gygi (56). 12 strong cation exchange (SCX) fractions were subjected to phosphopeptide enrichment using IMAC-Select Affinity Gel (Sigma-Aldrich) and subsequent peptide desalting with Stagetips (57). Samples were analyzed on an LTQ Orbitrap Velos mass spectrometer (Thermo Fisher Scientific) using a data-dependent Top10-MS2 method using (higher-energy collisional dissociation) HCD for reporter ion quantitation. Peptide identification and filtering was performed following the methods of Dephore and Gygi (56) but using a composite *Drosophila melanogaster* protein database. Data normalization and phosphosite localization was performed as previously described (58). The phosphorylation sites indicated above for Moira and MBD-R2 were localized with near certainty using the Ascore algorithm (Ascore > 19).

Protein identifier	Symbol	Phosphosite	Fold change Log2(10'/0')	Akt/RSK/S6 consensus motif (RxxS/T) (47)	Phosphosite
FBpp0082692	Mor	PGKRKRS#PAVVHK	0.30	Yes	Ser <sup>327</sup>
FBpp0082081	MBD-R2	KRASTGS#LGGSSG	0.26	No	Ser <sup>288</sup>