# Molecular grammars of intrinsically disordered regions that span the human proteome
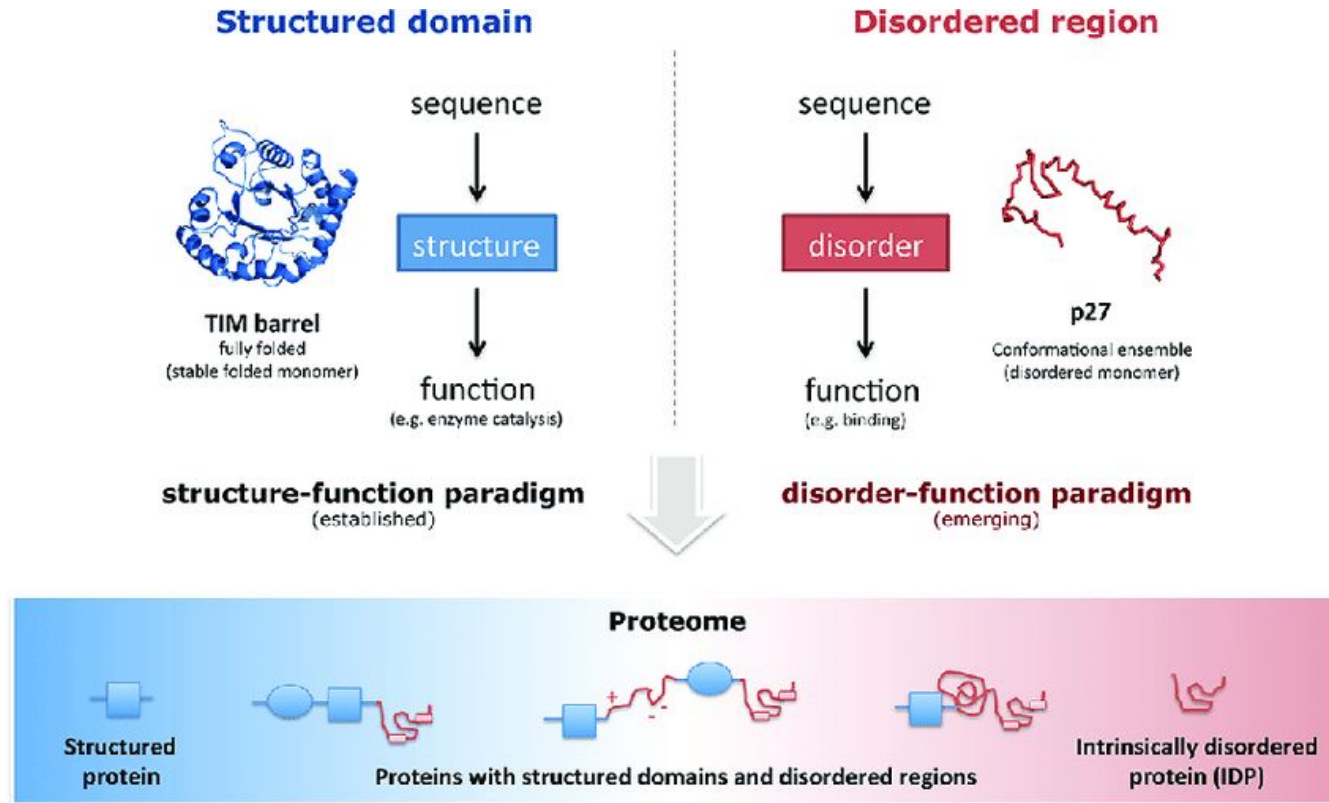
Kiersten M. Ruff, Matthew R. King, Alexander W. Ying Vicky Liu, Avnika Pant, Whitney E. Lieberman, Min Kyung Shinn,  Xiaolei Su, Cigall Kadoch, Rohit V. Pappu

LifeLU reading group

presented by Özdeniz Dolu

10.04.2025

# Intrinsically Disordered Region



2

# Intrinsically Disordered Region

- In human proteome, more than 50% of proteins contain at least one IDR
- IDRs show conformational heterogeneity and poor sequence conservation

Authors attempt to associate molecular grammars with IDRs:

- "IDR-specific molecular grammars are defined jointly by the non-random amino acid composition and the non-random patterning of distinct pairs of amino acid types with respect to one another."

# Proposed Approach

**NARDINI+ Algorithm**: Combination of NARDINI algorithm and compositional analyses inspired by other work in the literature.
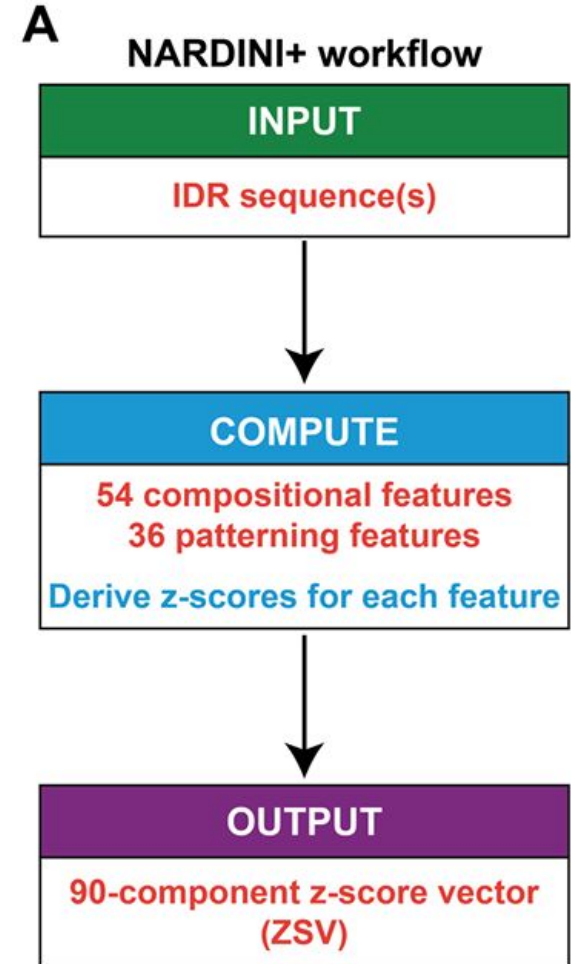
Through unsupervised learning methods and clustering, characterize IDRs in human proteome based on molecular grammars.

GIN: Grammars Inferred using NARDINI+

IDRome-spanning basis set, 30 clusters, unique fingerprint for each cluster

# NARDINI+

- 36 non-random binary patterns computed by NARDINI algorithm.
- Statistics on alphabet symbols used in the IDR yield 54 compositional features.
- Combined into a z-score vector.

**A**

**NARDINI+ workflow**

| INPUT |
| --- |
| IDR sequence(s) |

↓

| COMPUTE |
| --- |
| 54 compositional features<br>36 patterning features<br><br>Derive z-scores for each feature |

↓

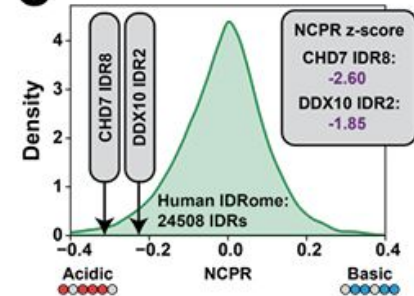| OUTPUT |
| --- |
| 90-component z-score vector<br>(ZSV) |

5

# NARDINI+ (Compositional Features)

- The fraction of each of the 20 naturally occurring residues (20 features)
- Fractions of polar, aliphatic, aromatic, (Lys + Arg), (Asp + Glu), charged residue, residues that promote chain expansion, disorder promoting residues (8 features)
- "20 values that quantify the presence of specific residue or RG patches"
- Arg/Lys ratio, Glu/Asp ratio
- Net charge per residue (NCPR) (4 features)
- Apparent isoelectric point (pI), Kyte-Doolittle hydrophobicity

**B**

CHD7 IDR8: VGSSEEKAADKAEGGPFKDGETLEGSDAEESLDKTAESSLLEDEIAQGEELDSLDGGDEIENNENDE

DDX10 IDR2: QKGGKRLEGTEHRQDNDTGNEEQEEEEDDEEEMEEKLAKAKGSQAPSLPNTSEAQKIKEVPTQFLDRDEEEEDAD

# NARDINI+ (Compositional Features)

- The fraction of each of the 20 naturally occurring residues (20 features)
- Fractions of polar, aliphatic, aromatic, (Lys + Arg), (Asp + Glu), charged residue, residues that promote chain expansion, disorder promoting residues (8 features)
- "20 values that quantify the presence of specific residue or RG patches"
- Arg/Lys ratio, Glu/Asp ratio
- Net charge per residue (NCPR) (4 features)
- Apparent isoelectric point (pI), Kyte-Doolittle hydrophobicity

**B**

CHD7 IDR8: VGSSEEKAADKAEGGPFKDGETLEGSDAEESLDKTAESSLLEDEIAQGEELDSLDGGDEIENNENDE

DDX10 IDR2: QKGGKRLEGTEHRQDNDTGNEEQEEEEDDEEEMEEKLAKAKGSQAPSLPNTSEAQKIKEVPTQFLDRDEEEEDAD

**C**

Density

NCPR z-score
CHD7 IDR8:
-2.60
DDX10 IDR2:
-1.85

CHD7 IDR8
DDX10 IDR2

Human IDRome:
24508 IDRs

-0.4    -0.2    0.0    0.2    0.4
Acidic          NCPR          Basic

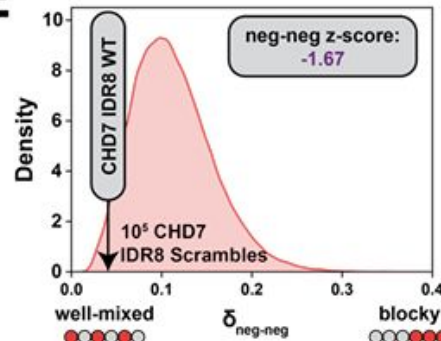# NARDINI+ (Binary Patterning Features)

- 36 binary features computed by NARDINI algorithm.
- Grouping of residues into 8 groups: polar, hydrophobic, positively charged, negatively charged, aromatic, alanine, glycine, and proline
- "For each unique pair of residue types U and X, the NARDINI algorithm computes a parameter δUX that quantifies the extent to which U and X are mixed or segregated along the linear sequence."
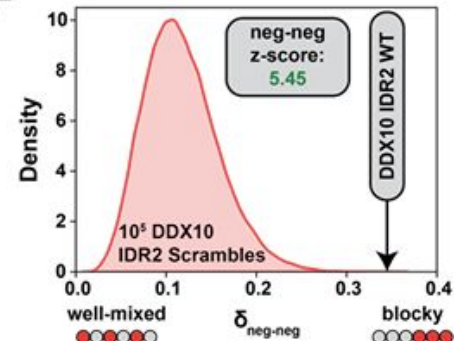


Patterning of residues of type U vs. X

- pol≡Polar: S, T, N, Q, C, H
- hyd≡Hydrophobic: I, L, V, M
- pos≡Positive: K, R
- neg≡Negative: D, E
- aro≡Aromatic: F, W, Y
- ala≡Alanine: A
- pro≡Proline: P
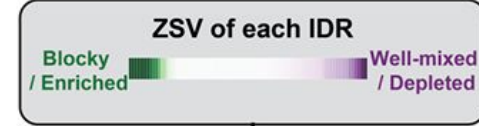- gly≡Glycine: G

8

# NARDINI+ (An Example)

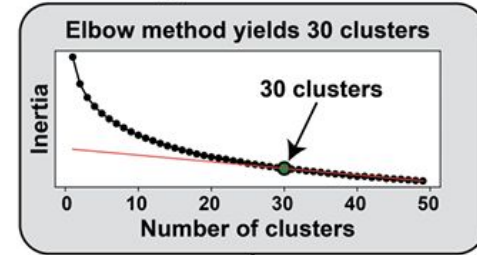Calculated features for 2 IDRs. Zero features are omitted.

# Grammars Inferred using NARDINI+ (GIN)

- 4,529 IDRs in preferred length range (out of 24508) used to cluster.
- K-means clustering with K = 30
- **All human IDRs** then mapped to clusters.
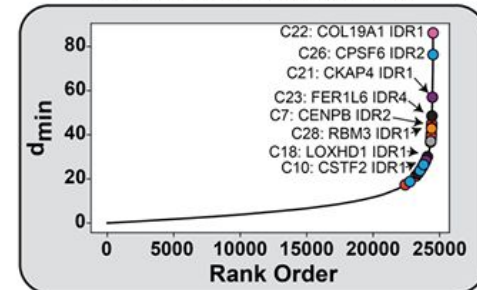- 30 clusters are referred as **GIN clusters**.



NARDINI+ applied to cluster all 24508 human IDRs

ZSV of each IDR

Blocky / Enriched — Well-mixed / Depleted

K-means clustering of IDRs of length $100 \leq n \leq 300$

Elbow method yields 30 clusters

30 clusters

Inertia

Number of clusters

Map all human IDRs onto 30 GIN clusters

C22: COL19A1 IDR1
C26: CPSF6 IDR2
C21: CKAP4 IDR1
C23: FER1L6 IDR4
C7: CENPB IDR2
C28: RBM3 IDR1
C18: LOXHD1 IDR1
C10: CSTF2 IDR1

$d_{min}$

Rank Order

# Grammars Inferred using NARDINI+ (GIN)



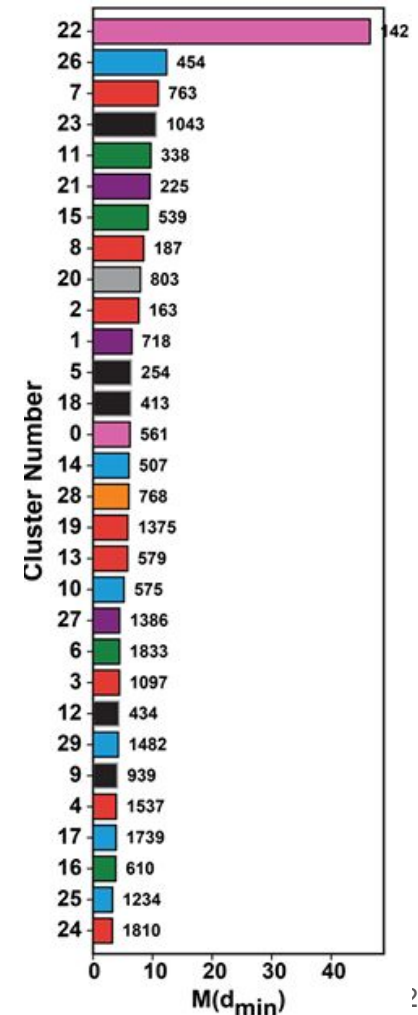Summary of the 30 GIN clusters than span the human IDRome

11

# Grammars Inferred using NARDINI+ (GIN)

Higher M(dmin) (Median of minimum inter-cluster distance) values imply stronger mapping to clusters.

Higher number of IDRs in a cluster should result in a weaker mapping. However there are some exceptions.

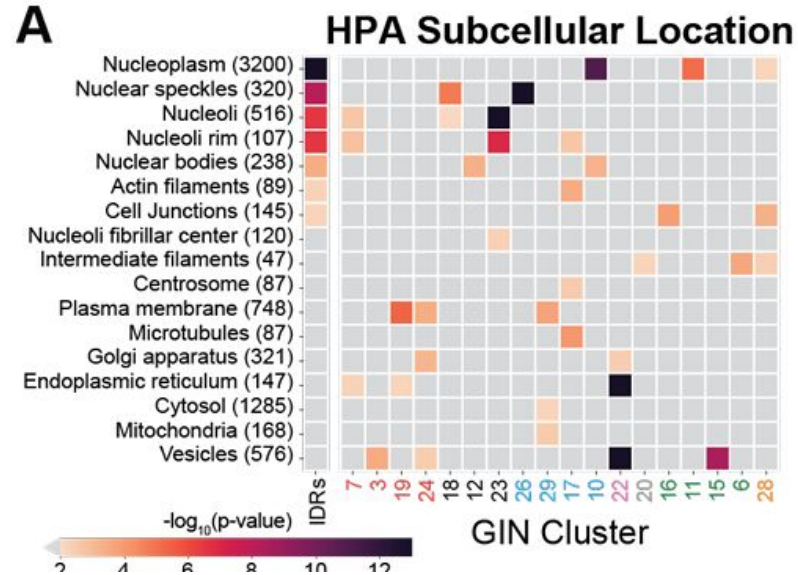For example, cluster 22 is very strongly mapped. 22 is uniquely defined by uniform distribution of Pro and Gly. Defining grammar for elastomers and collagens.



| Cluster Number | M(d_min) |
|---|---|
| 22 | 142 |
| 26 | 454 |
| 7 | 763 |
| 23 | 1043 |
| 11 | 338 |
| 21 | 225 |
| 15 | 539 |
| 8 | 187 |
| 20 | 803 |
| 2 | 163 |
| 1 | 718 |
| 5 | 254 |
| 18 | 413 |
| 0 | 561 |
| 14 | 507 |
| 28 | 768 |
| 19 | 1375 |
| 13 | 579 |
| 10 | 575 |
| 27 | 1386 |
| 6 | 1833 |
| 3 | 1097 |
| 12 | 434 |
| 29 | 1482 |
| 9 | 939 |
| 4 | 1537 |
| 17 | 1739 |
| 16 | 610 |
| 25 | 1234 |
| 24 | 1810 |

# Results: GIN as a resource

- Authors provide 2 Google Colab notebooks.
- First notebook takes gene id or uniprot accessions and calculates GIN annotations and ZSV vectors.
- Second notebook takes IDRs as fasta or a list of sequences and make the calculations from scratch.
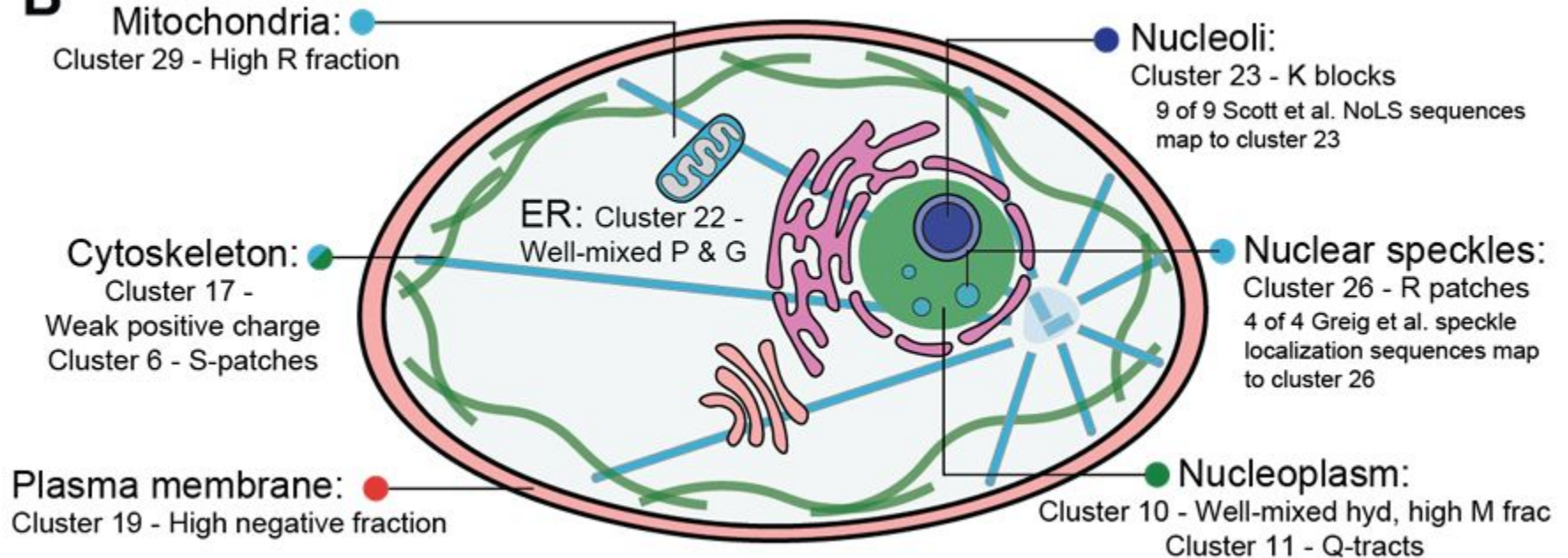
https://github.com/kierstenruff/RUFF_KING_Grammars_of_IDRs_using_NARDINI-

# Results: GIN clusters -> Sub-nuclear Localization

- Annotations of subcellular locations of proteins with IDRs from the Human Protein Atlas (HPA)
- Focus on IDRs of length ≥ 70 and non-linker IDRs of length ≥ 50 and high dmin values.
- Cluster 26 strongly associated with Nuclear speckles
- Cluster 23 strongly associated with Nucleoli rim

# Results: GIN clusters -> Sub-nuclear Localization

**B**

**Mitochondria:**
Cluster 29 - High R fraction

**Cytoskeleton:**
Cluster 17 -
Weak positive charge
Cluster 6 - S-patches

**Plasma membrane:**
Cluster 19 - High negative fraction

ER: Cluster 22 -
Well-mixed P & G

**Nucleoli:**
Cluster 23 - K blocks
9 of 9 Scott et al. NoLS sequences
map to cluster 23

**Nuclear speckles:**
Cluster 26 - R patches
4 of 4 Greig et al. speckle
localization sequences map
to cluster 26

**Nucleoplasm:**
Cluster 10 - Well-mixed hyd, high M frac
Cluster 11 - Q-tracts

15

# Results: GIN clusters -> Sub-nuclear Localization

For further testing the case for cluster 23 and 26:

- 5 monomeric proteins from HPA dataset whose IDRs are are clustered into either 23 or 26 are selected. These proteins have ambiguous annotations as to sub-nuclear localizations.
- Experiments are made on germinal vesicles (GVs) from live Xenopus laevis oocytes.
- Germinal Vesicle: Nucleus of an oocyte
- Oocyte: A developing egg

# Results: GIN clusters -> Sub-nuclear Localization

- Except for GFP, proteins with cluster 23 are partitioned into Nucleolus and proteins with cluster 26 are partitioned into Nuclear Speckles.

- The **partition coefficient**, abbreviated **P**, is defined as a particular ratio of the concentrations of a solute between the two solvents (a biphase of liquid phases) (Wikipedia)

# Results: GIN clusters -> Sub-nuclear Localization

Two naturally occurring proteins GPatch3 and GPatch4 are investigated. Both contain GPatch domain. However, their localization preference are different.

Localization preferences of proteins same domains can be determined by IDRs.

Authors note that, in some cases the vice versa is true.

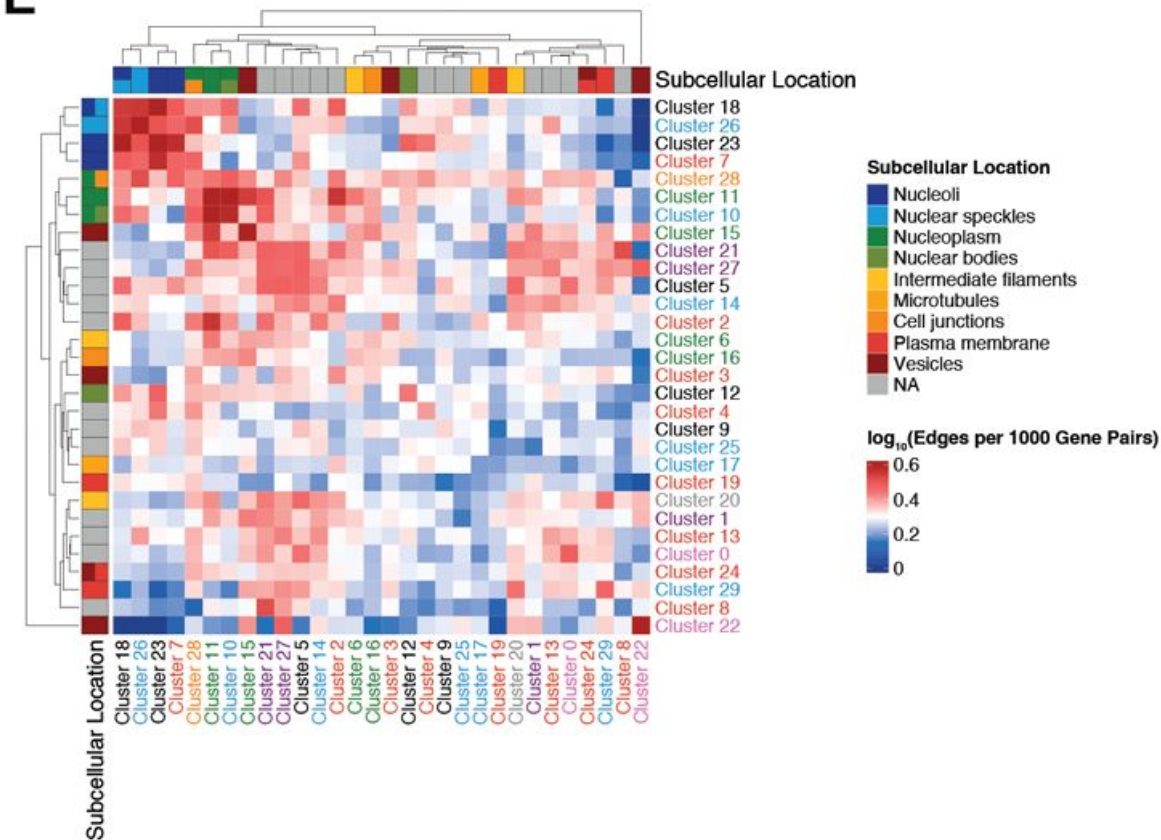# Results: GIN clusters -> MF and BP

# Results: GIN clusters -> MF and BP

**DepMap24Q4** Dataset

- "Genome-wide CRISPR knockout screening in over 1000 cancer cell lines"
- "Two genes (proteins) might be functionally linked if their fitness effects upon knockout across the cell lines are correlated"
- How proteins with IDR GIN cluster annotations are associated?

# Results: GIN clusters -> MF and BP

Heatmap quantifying functional relationships based on **DepMap24Q4** dataset.

Higher value corresponds to more relationship.

# Results: GIN clusters -> MF and BP

Inter and intra cluster fitness correlations based on **DepMap24Q4**.

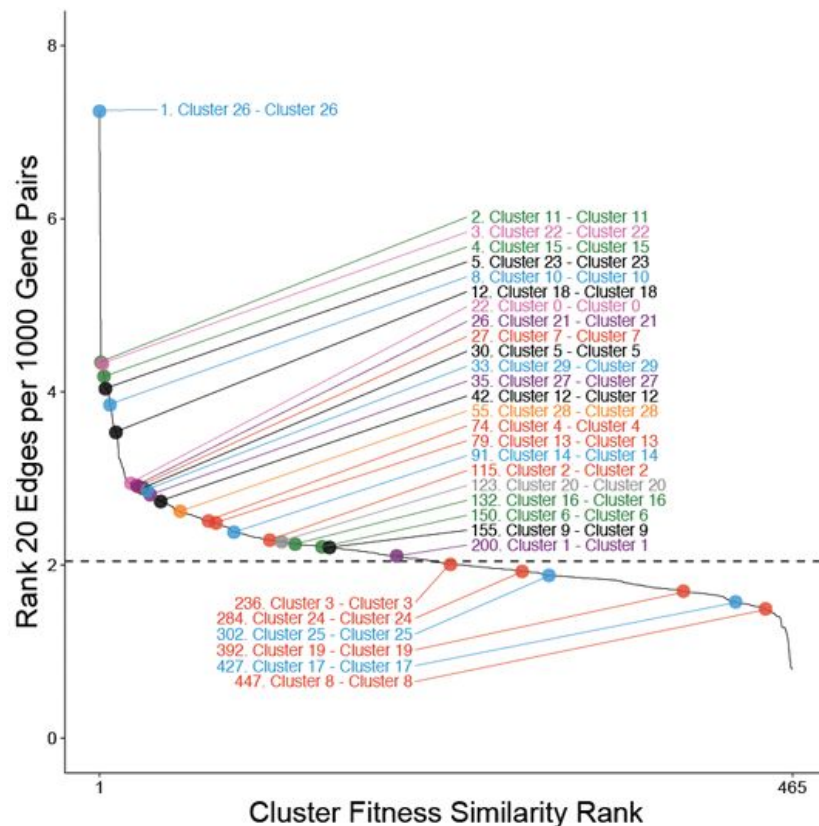Only the intra-cluster data shown on right.

An interesting note:

Top ranking clusters 26, 11, 22, 15, 23, and 10 are also highly enriched for specific subcellular locations. Suggesting that subcellular location might be an important signal for fitness data.

# Results: GIN clusters -> MF and BP

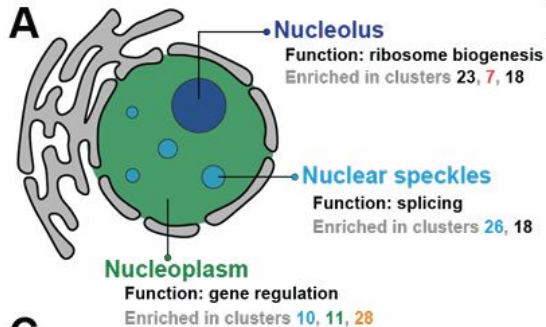6 of 10 top-ranked gene (protein) fitness correlations are (within GIN clusters)

Genes/proteins that are associated some GIN clusters seem to have functional relationships.

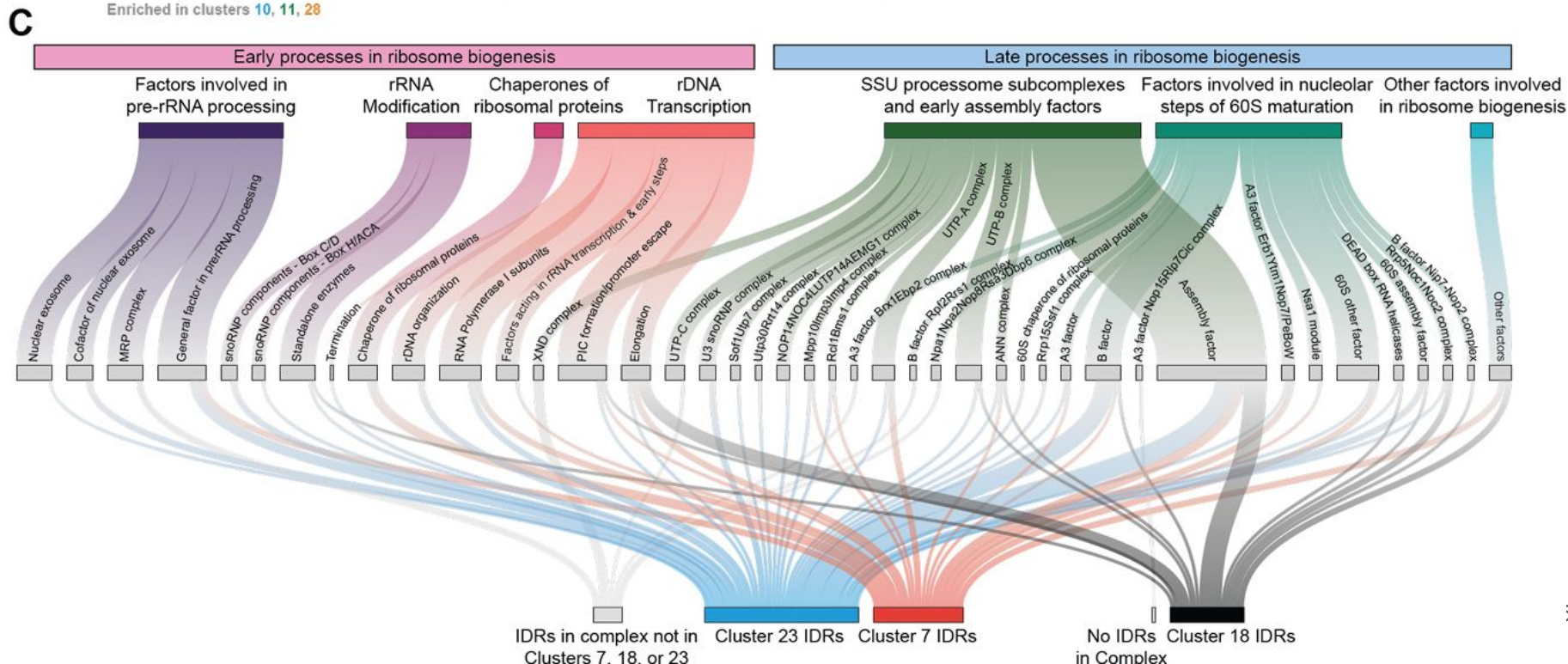# Results: GIN clusters -> Discrete Processes

Closer look into by combining known functions of subcellular locations and GIN cluster enrichment.
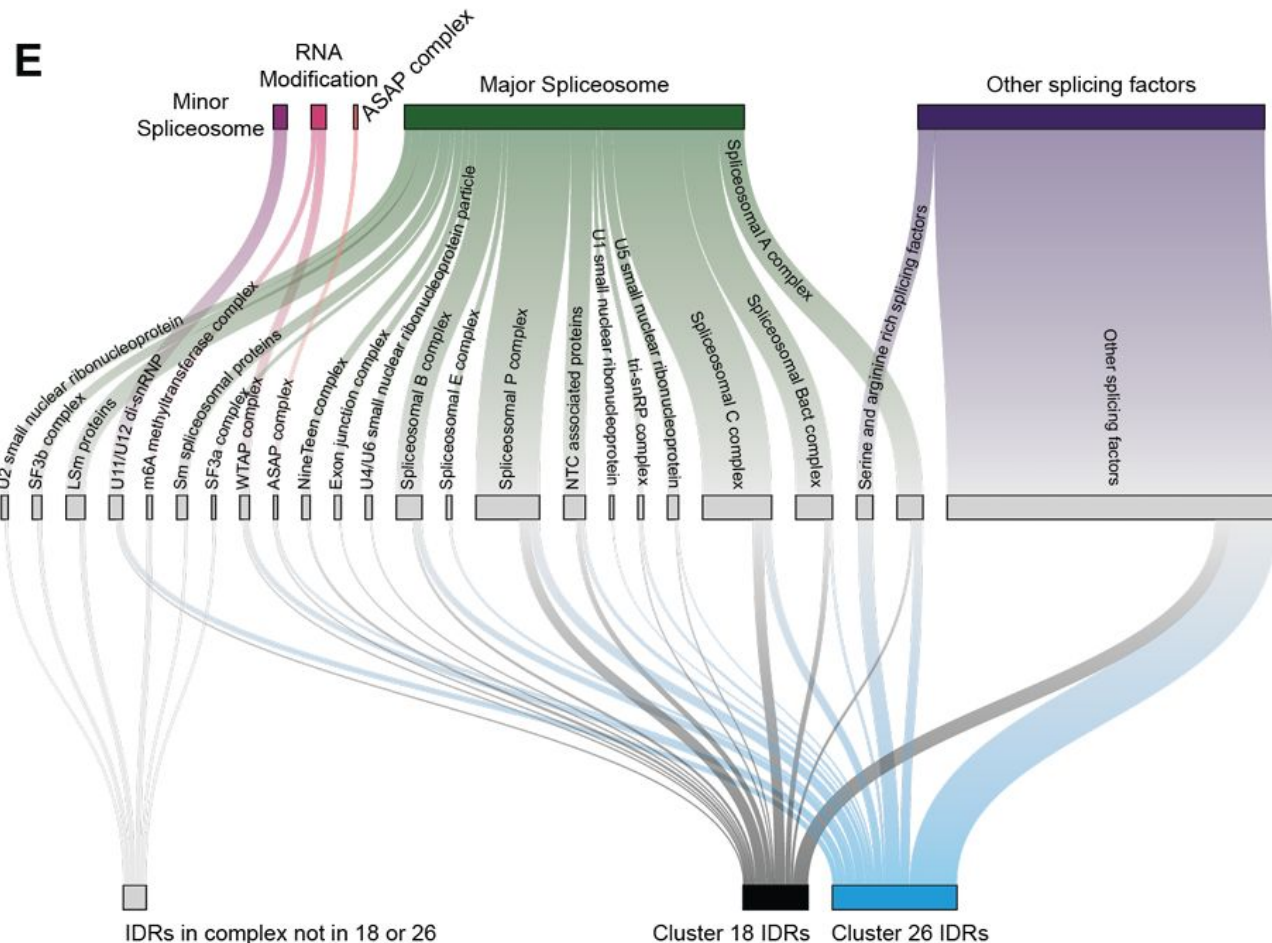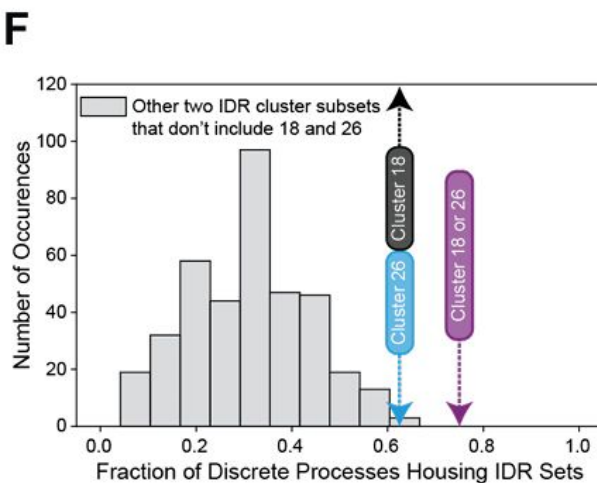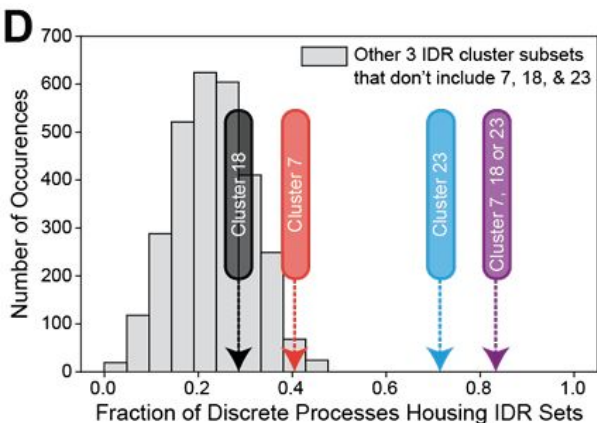
**A**

Nucleolus
**Function: ribosome biogenesis**
Enriched in clusters 23, 7, 18

Nuclear speckles
**Function: splicing**
Enriched in clusters 26, 18

Nucleoplasm
**Function: gene regulation**
Enriched in clusters 10, 11, 28

**B**

| | GIN | IDR Grammar | Example IDR | Example Sequence |
|---|---|---|---|---|
| Nucleoli | 23 | K blocks | POLR1F IDR1 | ...KKKKKKKKHQEVQDQDPVFQGSDSSGYQSDHKKKKKRKHSEEAEFTPPL... |
| | 7 | D/E-tracts | UBTF IDR4 | ...EEDDEEDEDDDEDEDEEEDDENGDSSEDGGDSSESSSEDESEDGDENEED... |
| | 18 | Large negative blocks with positive blocks | NCL IDR1 | ...DDDDDEEDDSEEEAMETTPAKGKKAAKVVPVKAKNVAEDEDEEEDDEDED... |
| NS | 26 | R patches | SRRM2 IDR2 | ...RGRSRSRTPARRGRSRSRTPARRRSRSRTPTRRRSRSRTPARRGRSRSRT... |
| | 18 | Large negative blocks with positive blocks | CD2BP2 IDR1 | ...DEEDEDEIIVPKKKLVDPVAGSGGPGSRFKGKHSLDEEEEDDDDGSSGK... |
| NP | 10 | Well-mixed hydrophobics, enriched in M | ARID1B IDR4 | ...MSMPDVMGRMPYPEPNKDPFGGMRKVPGSSEPFMTQGQMPNSSMQDMYNQS... |
| | 11 | Q-tracts | PAXIP1 IDR2 | ...QQFHQLQQHQLQQQQLAQLQQQHSLLQQQQQQIQQQQLQRMHQQQQQQQO... |
| | 28 | High aromatic fraction, specifically Ys | POLR2A IDR1 | ...YSPTSPSYSPTSPSYSPTSPSYSPTSPSYSPTSPSYSPTSPSYSPTSPSY... |

**C**

Early processes in ribosome biogenesis

Late processes in ribosome biogenesis

Factors involved in pre-rRNA processing · rRNA Modification · Chaperones of ribosomal proteins · rDNA Transcription

SSU processome subcomplexes and early assembly factors · Factors involved in nucleolar steps of 60S maturation · Other factors involved in ribosome biogenesis

IDRs in complex not in Clusters 7, 18, or 23 · Cluster 23 IDRs · Cluster 7 IDRs · No IDRs in Complex · Cluster 18 IDRs

**D** Number of Occurrences vs Fraction of Discrete Processes Housing IDR Sets. Legend: Other 3 IDR cluster subsets that don't include 7, 18, & 23. Labels: Cluster 18, Cluster 7, Cluster 23, Cluster 7, 18 or 23.

**E** Minor Spliceosome, RNA Modification, ASAP complex, Major Spliceosome, Other splicing factors. Components: U2 small nuclear ribonucleoprotein, SF3b complex, LSm proteins, U11/U12 di-snRNP, m6A methyltransferase complex, Sm spliceosomal proteins, SF3a complex, WTAP complex, ASAP complex, NineTeen complex, Exon junction complex, U4/U6 small nuclear ribonucleoprotein, Spliceosomal B complex, Spliceosomal E complex, Spliceosomal P complex, NTC associated proteins, U1 small nuclear ribonucleoprotein, U5 small nuclear ribonucleoprotein, tri-snRP complex, Spliceosomal C complex, Spliceosomal A complex, Spliceosomal Bact complex, Serine and arginine rich splicing factors, Other splicing factors. IDRs in complex not in 18 or 26, Cluster 18 IDRs, Cluster 26 IDRs.

**F** Number of Occurrences vs Fraction of Discrete Processes Housing IDR Sets. Legend: Other two IDR cluster subsets that don't include 18 and 26. Labels: Cluster 18, Cluster 26, Cluster 18 or 26.

26

# Results: Exceptional Grammars

All 24508 IDRs in human proteome is sorted by a "given" GIN cluster feature.

IDRs in 99th percentile are examined.

# Results: Exceptional Grammars



**The earliest ribosome biogenesis processes are enriched in top scoring K blocks IDRs in the human IDRome (24,508 IDRs)**

# Results: Exceptional Grammars



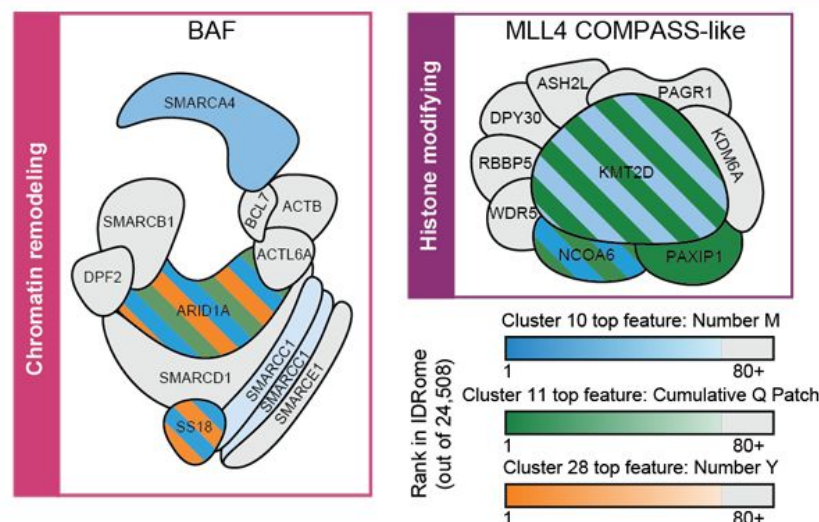The earliest spliceosome processes are enriched in top scoring R patch IDRs in the human IDRome

# Results: Exceptional Grammars



Core enhancer activation complexes are enriched in top scoring nucleoplasmic features in the human IDRome

# Results: RNA Pol I & II and evolutionarily conserved exceptional grammars

RNA polymerases play an essential role in transcription.

The subunit-specific IDRs of Pol I and Pol II have grammars with exceptional features that define clusters 23 and 28, respectively.

Orthologs of these IDRs throughout IDRomes of 8 species is investigated.

# Results: RNA Pol I & II and evolutionarily conserved exceptional grammars

# Results: RNA Pol I & II and evolutionarily conserved exceptional grammars

Although sequences differ throughout,, "exceptionality" of POLR1F is conserved in many species.



**B** Exceptional grammar features of Pol I and Pol II specific IDRs

**C** POLR1F: pos-pos z-score

# Results: RNA Pol I & II and evolutionarily conserved exceptional grammars

Figures and analyses containing further investigation and atomic simulations are also included in Figure 6. However, omitted from this presentation.

# Results: Cancer mutations disrupt exceptional grammars

Cancer driver genes are investigated. Proteins associated with them are seen to be enriched in clusters 10, 11, 12.

# Results: Cancer mutations disrupt exceptional grammars

Cancer driving mutations of C-terminus IDRs of CREBBP and EP300 can lead to "exceptionality loss" for grammars.

# Results: Cancer mutations disrupt exceptional grammars



Gain / loss of exceptional IDR grammars in patient derived fusions with DNA binding domains

# Thanks for listening