# From Sequences to Molecules: An open-source Monomer-Centric Toolkit

**Davit Rizhinashvili, Datagrok**
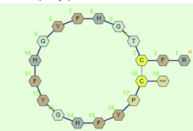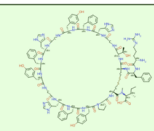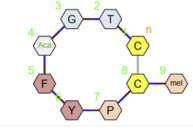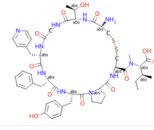
**RDKit UGM 2025, Prague**

**September 11, 2025**

# Bridging sequences and molecules - Why this is valuable

- **Sequence representations are more human readable**

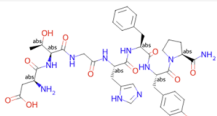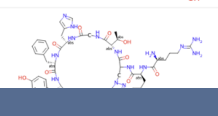- **Modular building block approach**

- **Rapid sequence modification**

- **Efficient database searches**

- **Easier computational analysis, access to sequence based tools (MSA, composition analysis, clustering, etc.)**

- **Error prevention**

- **Interactive UI**

**While also being able to calculate properties and conduct analysis on molecular level**

# Datagrok: enterprise-ready life sciences platform

- Data access, exploratory data analysis, scientific computing, etc

- Analyzing big datasets completely in the browser

- Proprietary core, open-source plugin ecosystem

- Industry adoption

- Domain-agnostic

- Cheminformatics as a plugin

- RDKit At its cheminformatics core
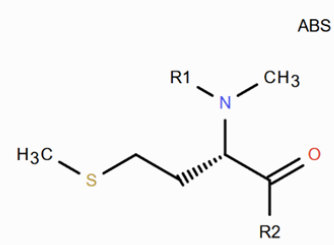
- Free for personal or academic use

# Monomer libraries



- **Manage multiple libraries**

- **Add, edit, delete monomers**

- **Selectively extract**

- **Manage duplicates**

- **Standardization and validation**

# Monomer Manager

# Monomers

meI   hHis   Aca   N   T dK   Thr…   Aca   D…   Tyr…   dV   Thr…   N

mel — N-Methyl-Isoleucine

HELMCore Library

dD   dK   Thr…   Aca   D…   Tyr…   dV   Trp_Me   N

hPh   dK   Thr…   Aca   D…   Tyr…   dV   Phe…   N

dI   hHis   Aca   N   T dK   Thr…   Aca   D…   Tyr…   dV   Thz   N

# Linear Sequences
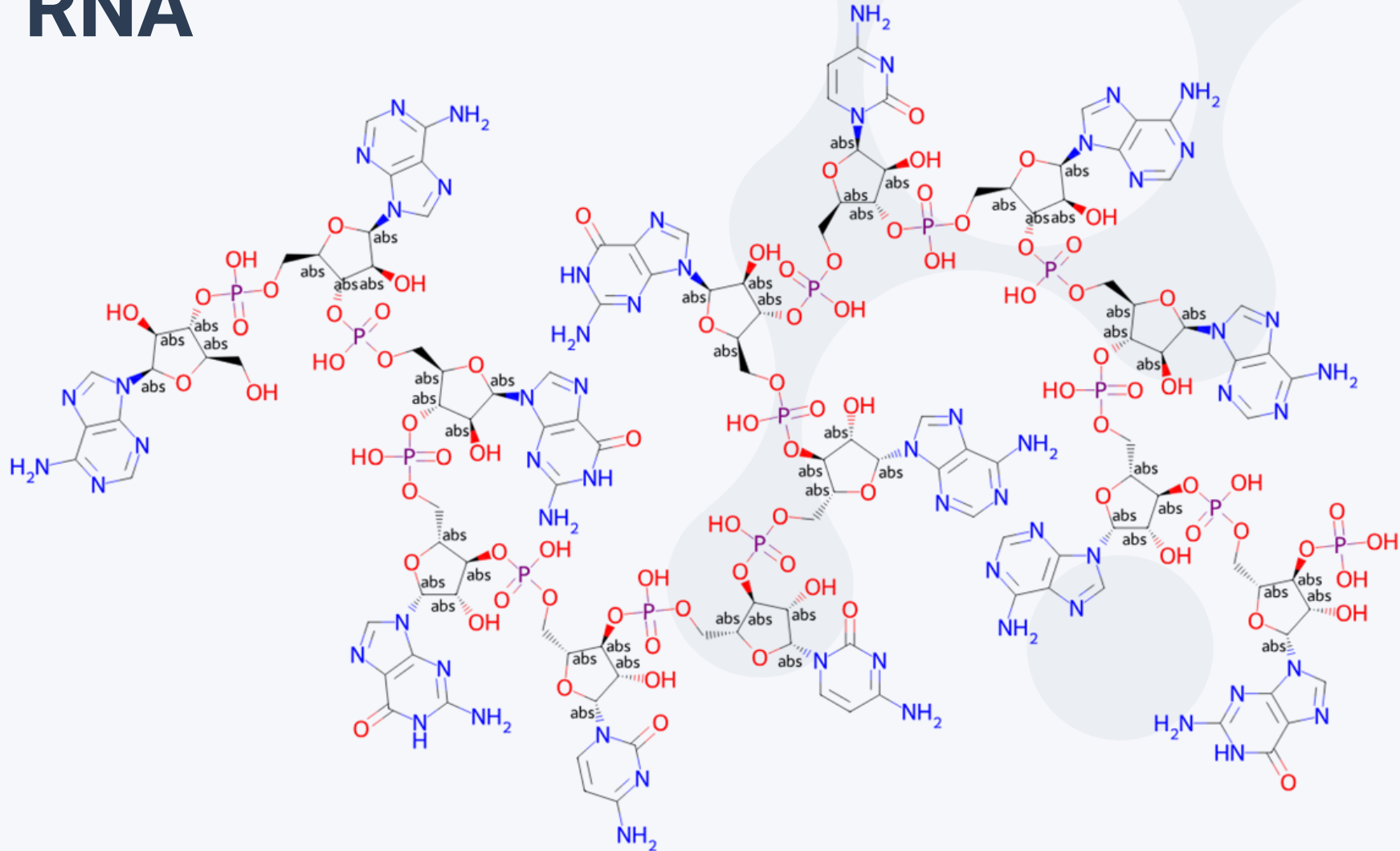
- R2 of First monomer connects to R1of second monomer
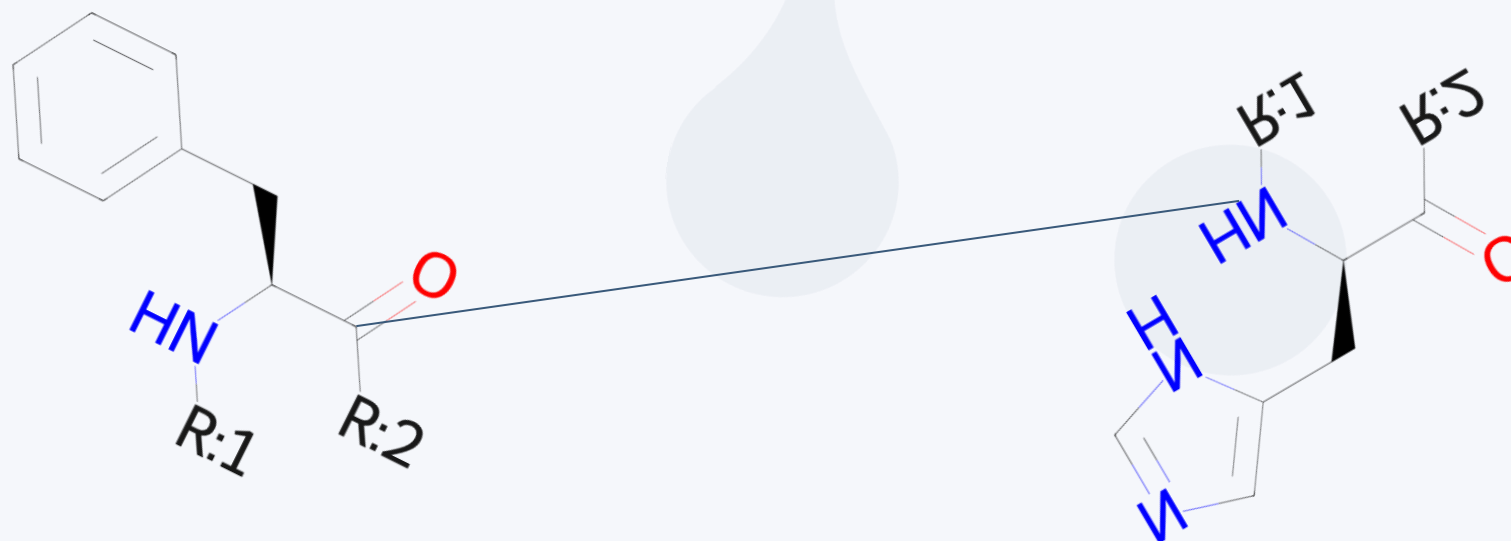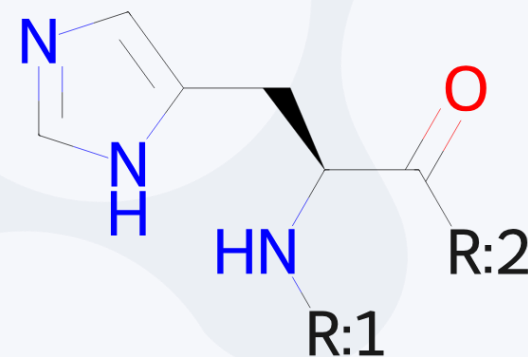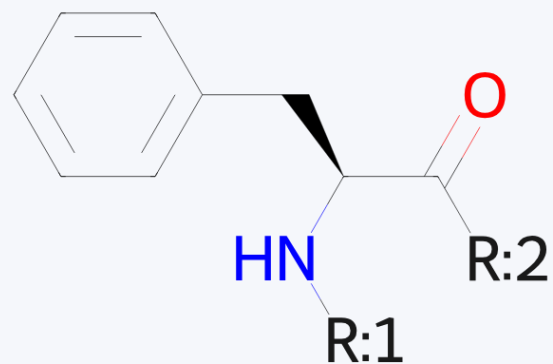
# Linear Sequences
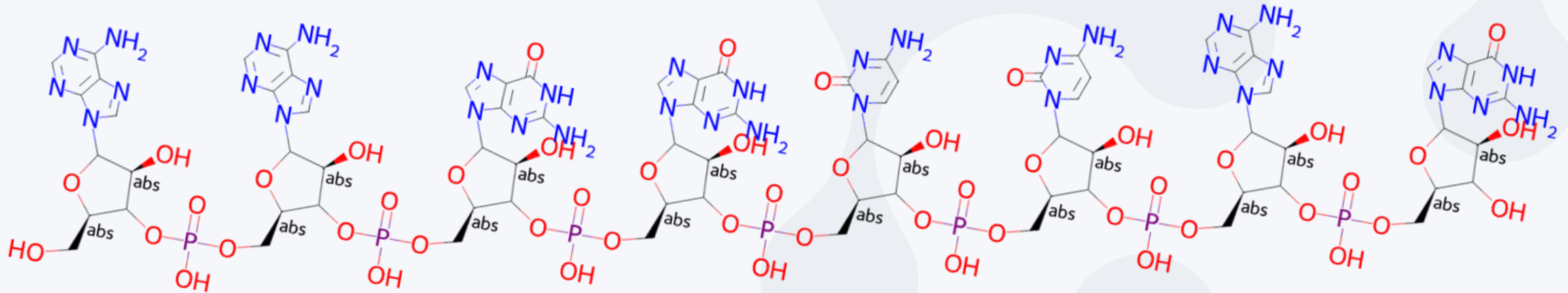## Peptides

# Linear Sequences
## DNA / RNA

# Linear Sequences
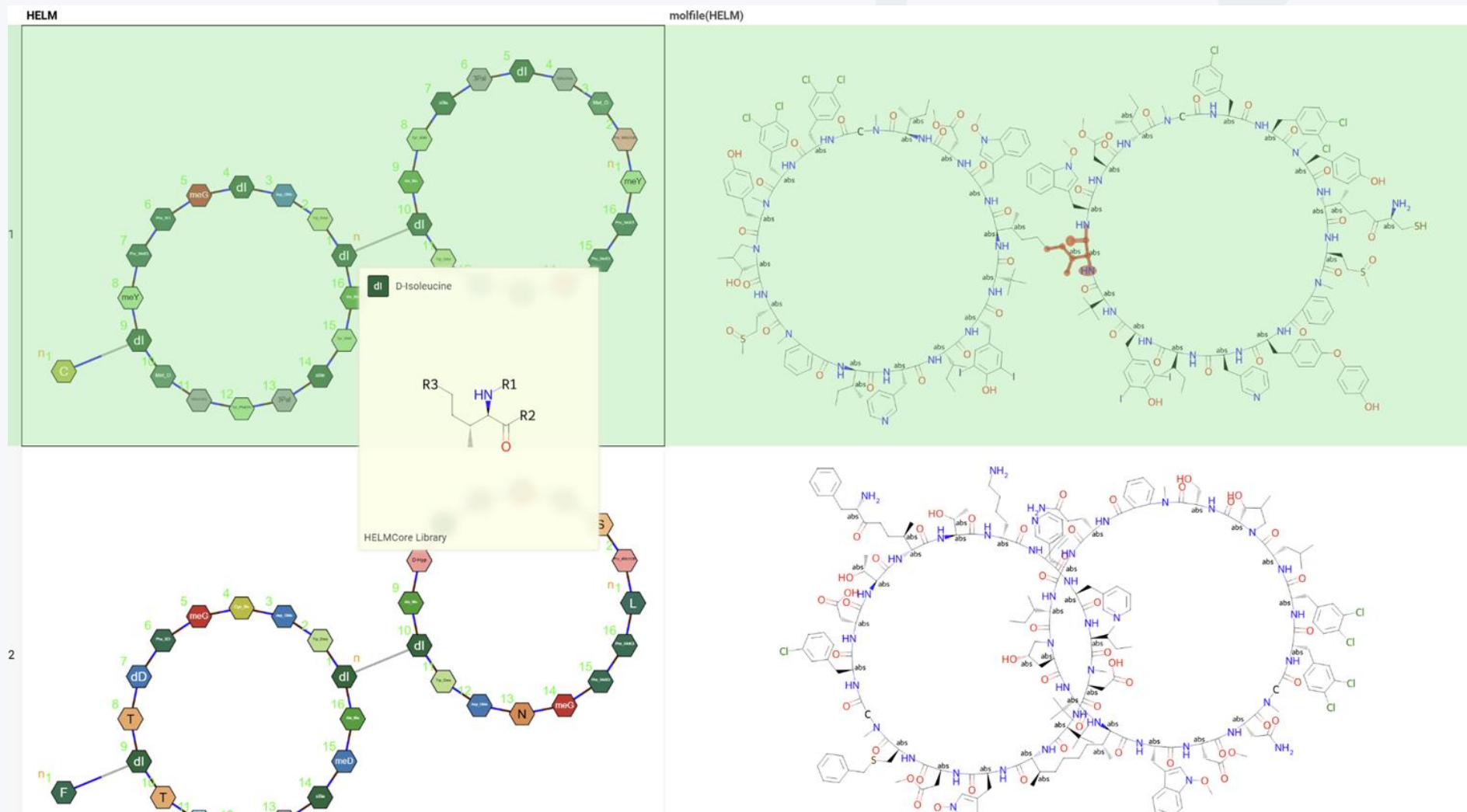## More human-friendly depiction
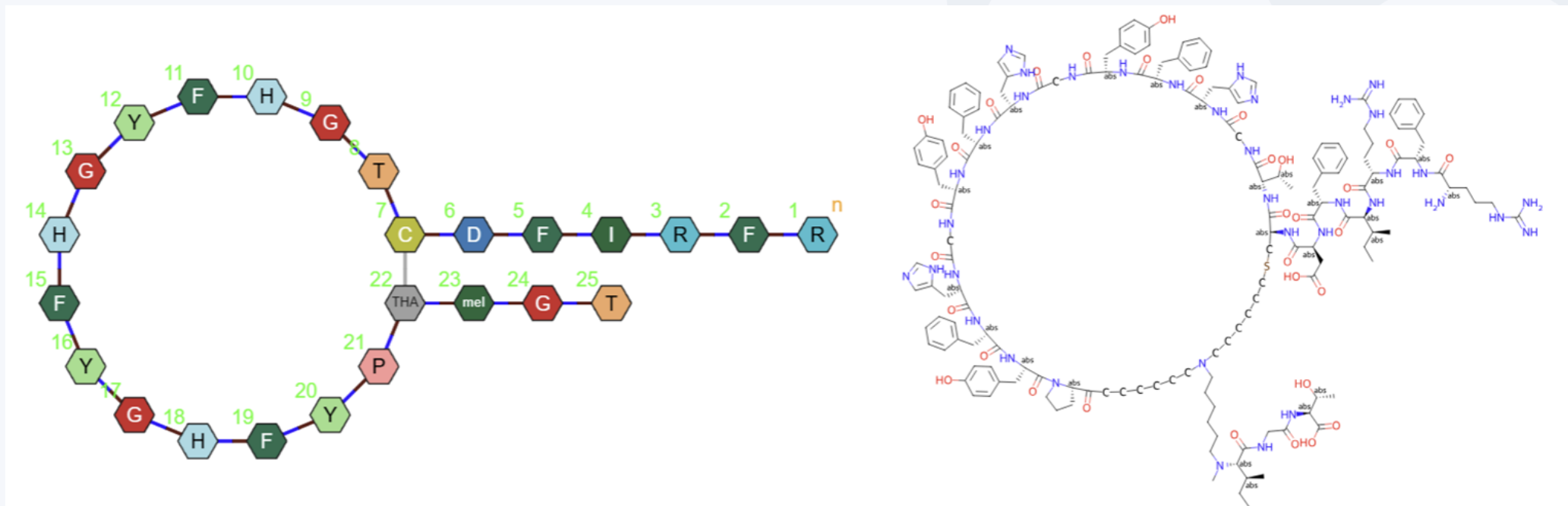
# Linear Sequences
## Peptides

# Linear Sequences
## DNA / RNA

# Cyclic structures: HELM

# HELM



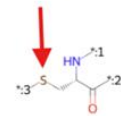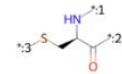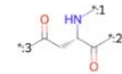PEPTIDE1{R.F.R.I.F.D.C.T.G.H.F.Y.G.H.F.Y.G.H.F.Y.P.[THA].[meI].G.T}$PEPTIDE1,PEPTIDE1,7:R3-22:R3$$$V2.0

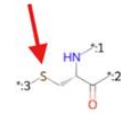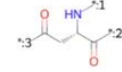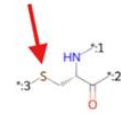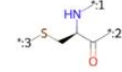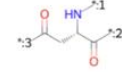# Rules - make sequences readable again

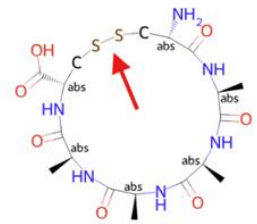PEPTIDE1{R.F.R.I.F.D.C.T.G.H.F.Y.G.H.F.Y.G.H.F.Y.P.[THA].[meI].G.T}$PEPTIDE1,PEPTIDE1,7:R3-22:R3$$$V2.0

R-F-R-I-F-D-C(1)-T-G-H-F-Y-G-H-F-Y-G-H-F-Y-P-THA(1)-meI-G-T

# Link Rules

# BILN - How is it different?

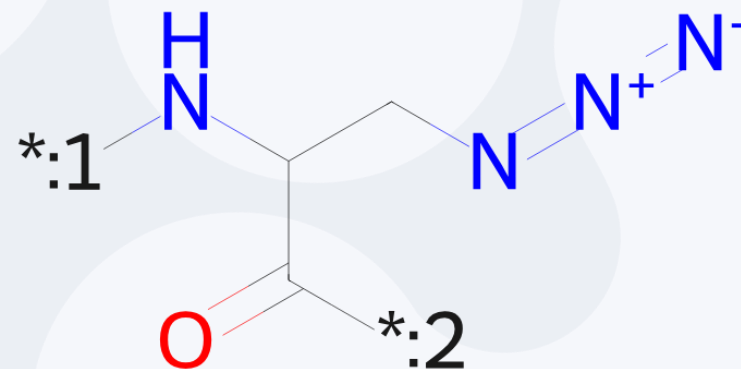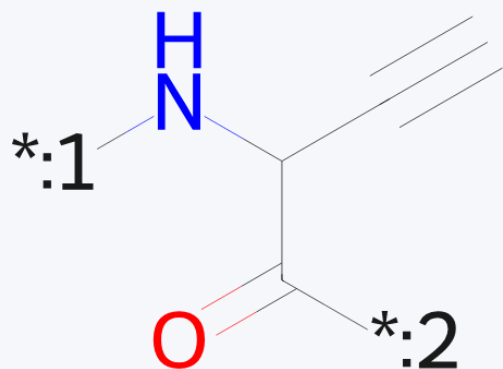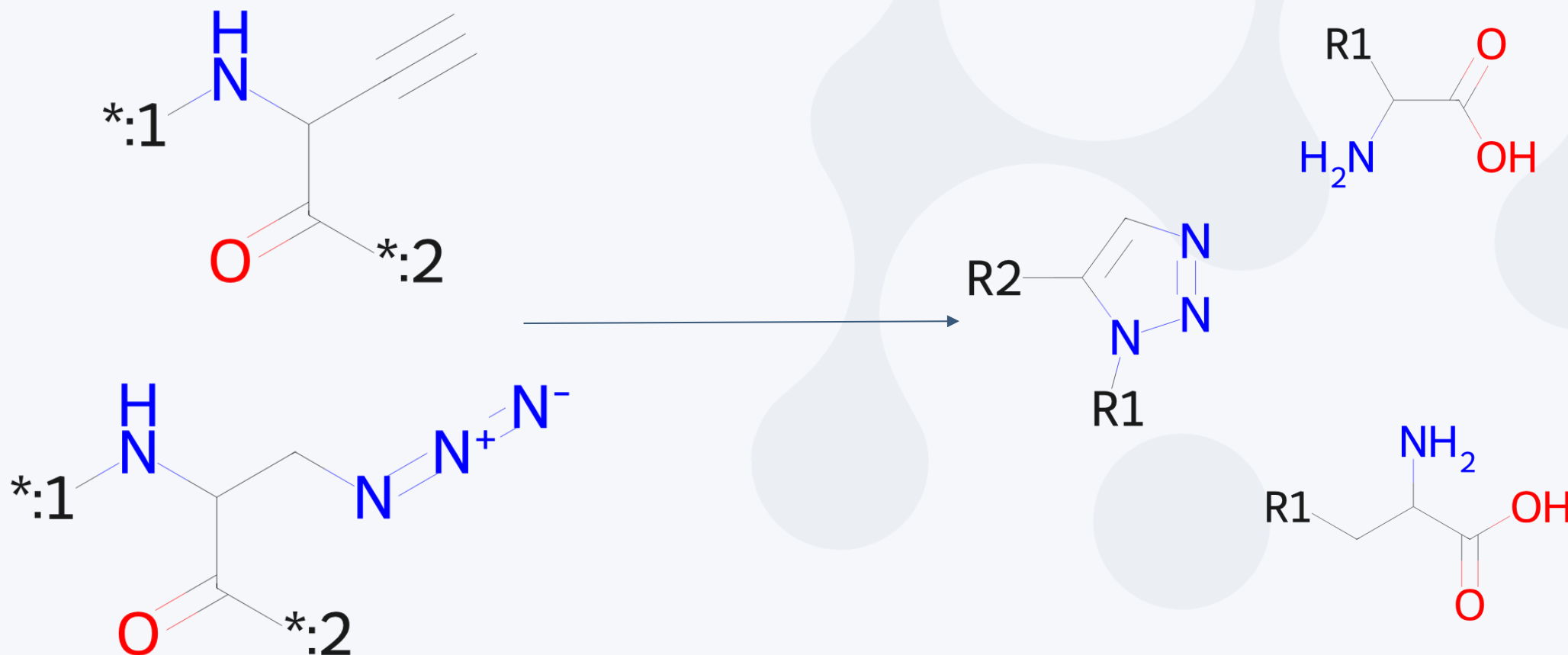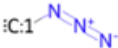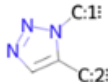| | |
|---|---|
| **BILN** | D-T-H-F-P-I-C(1,3)-I-F-C(2,3)-C(3,3)-G-C(2,3)-C(4,3)-H-R-S-K-C(3,3)-G-M-C(4,3)-C(1,3)-K-T |
| **HELM** | PEPTIDE1{D.T.H.F.P.I.C.I.F.C.C.G.C.C.H.R.S.K.C.G.M.C.C.K.T}$PEPTIDE1,PEPTIDE1,7:R3-23:R3\|PEPTIDE1,PEPTIDE1,10:R3-13:R3\|PEPTIDE1,PEPTIDE1,11:R3-19:R3\|PEPTIDE1,PEPTIDE1,14:R3-22:R3$$$ |
| **Structure** | D-T-H-F-P-I-C-I-F-C-C-G-C-C-H-R-S-K-C-G-M-C-C-K-T |

Fox, Thomas, et al. "BILN: a human-readable line notation for complex peptides." Journal of Chemical Information and Modeling 62.17 (2022): 3942-3947.

17

# Reaction Rules

# Reaction Rules

# Reaction Rules

| Rules | | | | | Examples | | |
|-------|--|--|--|--|----------|--|--|
| Name | First reactant | Second reactant | Product | Code | Monomers | Helm | molfile(sequence) |
| 1 GGaz | | | | 4 | azG(4)  A  A  A  A  aG(4) | | |
| 2 DARR | | | | 8 | DRR1(8)  A  A  A  A  DRR2(8) | | |
| 3 ODARR | | | | 8 | ODAR1(8)  A  A  A  A  ODAR2(8) | | |
| 4 PHAR_CIS | | | | 7 | PHAR1(7)  A  A  A  A  PHAR2(7) | | |
| 5 PHAR_TRANS | | | | 6 | PHAR1(6)  A  A  A  A  PHAR2(6) | | |

# Reaction Rules
## Automatically created reaction product monomers

# Conversion with rules

# Enumerator

# Enumerator

# Thank You!

## Acknowledgements

- RDKit maintainers and community

- Datagrok team

- Rhitankar Pal

- All our users ☺

datagrok