

Cheminformatics Meets Biology: Hybrid Molecular Representations in Drug Discovery



Jonathan Bisson
Alex Clark – Tad Hurst

IF YOU HAVE TO REMEMBER ONE THING



*Don't **ever** throw away the original file!*

Dr. Alex M. Clark



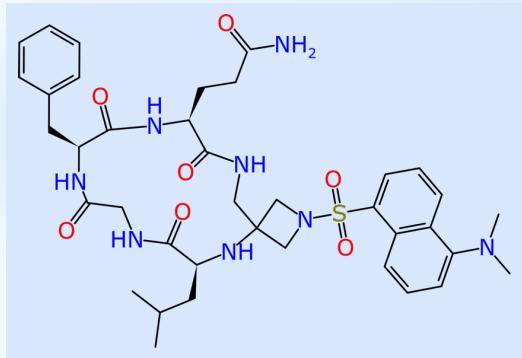
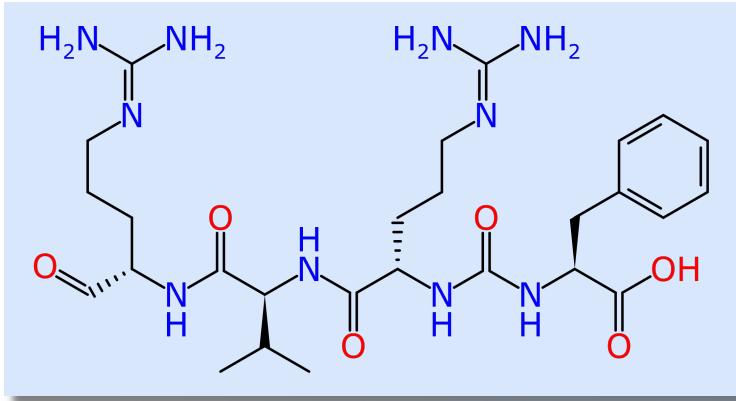
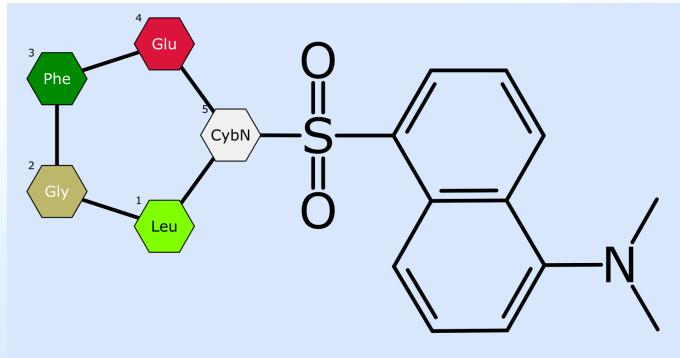
<https://cheminf20.org/2014/05/05/on-the-myth-of-chemical-structure-format-conversion/>

H	R	R	T	S
Q	R	C	E	E
E	R	Y	V	P
V	N	Y	C	C
Q	W	G	M	

Bioinformatics or Cheminformatics

We need both...

L,G,F,E,CybN-SO2NphNMe2,*



V3000 molfile feature: **SCSR**

Self-Contained Structure Representations

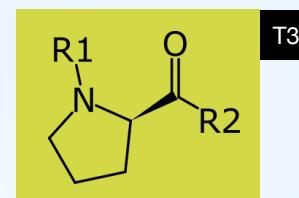
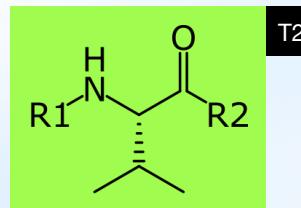
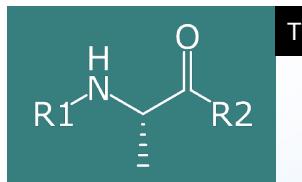
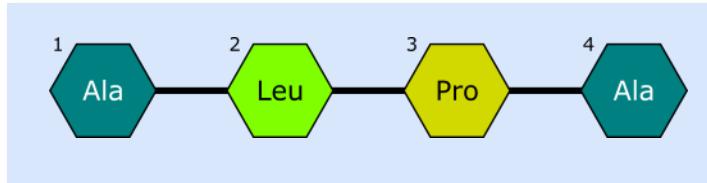
Layout information as 2D sketch

Monomers as **templates**: scalable

Self contained

Full **atomic** structure is implied

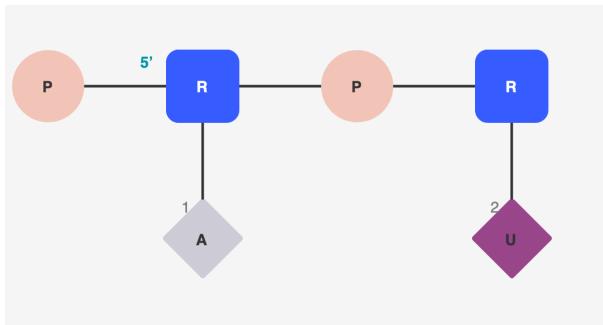
Contains the **sequence** data



https://www.rdkit.org/docs/RDKit_Book.html#id32

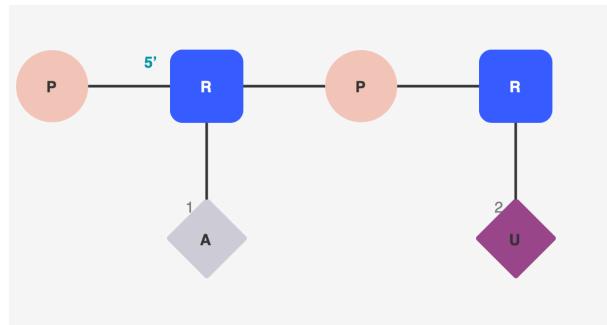
Generated by WebMolKit

```
0 0 0 0 0          999 V3000
M V30 BEGIN CTAB
M V30 COUNTS 6 5 0 0 1
M V30 BEGIN ATOM
M V30 1 P 0.0000 0.0000 0.0000 0 CLASS=PHOSPHATE SEQID=1 ATTCHORD=(2 3 Br)
M V30 2 P 3.0000 0.0000 0.0000 0 CLASS=PHOSPHATE SEQID=2 ATTCHORD=(4 3 Al 4 Br)
M V30 3 R 1.5000 0.0000 0.0000 0 CLASS=SUGAR SEQID=1 ATTCHORD=(6 1 Al 2 Br 5 Cx)
M V30 4 R 4.5000 0.0000 0.0000 0 CLASS=SUGAR SEQID=2 ATTCHORD=(4 2 Al 6 Cx)
M V30 5 A 1.5000 -1.5000 0.0000 0 CLASS=BASE SEQID=1 ATTCHORD=(2 3 Al)
M V30 6 U 4.5000 -1.5000 0.0000 0 CLASS=BASE SEQID=2 ATTCHORD=(2 4 Al)
M V30 END ATOM
```



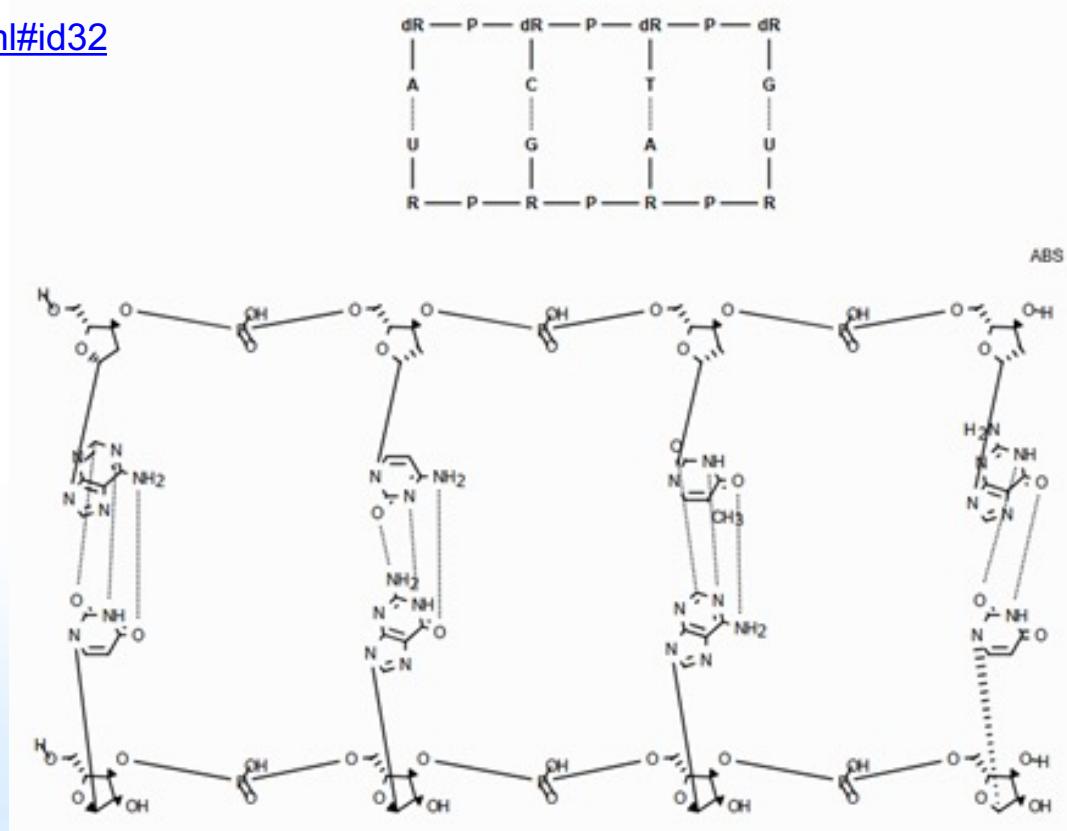
SCSR – Template and leaving groups

```
M V30 BEGIN TEMPLATE
M V30 TEMPLATE 1 PHOSPHATE/P/P/ NATREPLACE=PHOSPHATE/P
M V30 BEGIN CTAB
M V30 COUNTS 5 4 3 0 1
M V30 BEGIN ATOM
M V30 1 P 0.0000 0.0000 0.0000 0
M V30 2 O 1.5000 0.0000 0.0000 0
M V30 3 O -1.5000 0.0000 0.0000 0
M V30 4 O 0.0000 -1.5000 0.0000 0
M V30 5 O 0.0000 1.5000 0.0000 0
M V30 END ATOM
M V30 BEGIN BOND
M V30 1 1 1 2
M V30 2 1 1 3
M V30 3 1 1 4
M V30 4 2 1 5
M V30 END BOND
M V30 BEGIN SGROUP
M V30 1 SUP 0 LABEL=P ATOMS=(3 1 4 5) XBONDS=(2 2 1) CLASS=PHOSPHATE ... SAP=(3 1 3 Al) SAP=(3 1 2 Br)
M V30 2 SUP 0 LABEL=O ATOMS=(1 3) XBONDS=(1 2) CLASS=LGRP
M V30 3 SUP 0 LABEL=O ATOMS=(1 2) XBONDS=(1 1) CLASS=LGRP
M V30 END SGROUP
M V30 END CTAB
```



SCSR – Wobble bonds

https://www.rdkit.org/docs/RDKit_Book.html#id32



Content creation



Draw anything with EPAM'Ketcher

Bulk import: markup sequences

Custom monomers

CDD.VAULT • McKerrow Vault ▾

AI Chat · Help · Log Out

Explore Data

ELN

Inventory

Import Data

Reports

Settings

✉ 4

Full-Access User ▾

Step 1: Choose Data File and Parse

Step 2: Map Fields

Step 3: Commit Data

File: oligomers.csv

Project: McKerrow Vault · Owner: Full-Access User

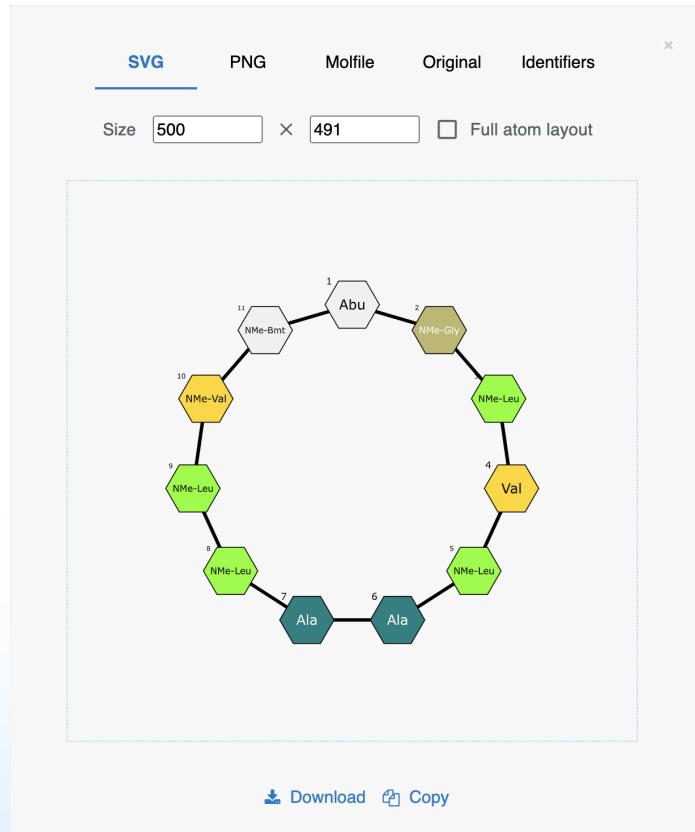
Compose macromolecules from columns [Edit](#)

Not part of macromolecule Wrapping: Off On
 Linear peptide sequence Width: 5 units
 Cyclic peptide sequence
 Nucleotide sequence
 Antibody chain
 Translate DNA/RNA codons into amino acid codes

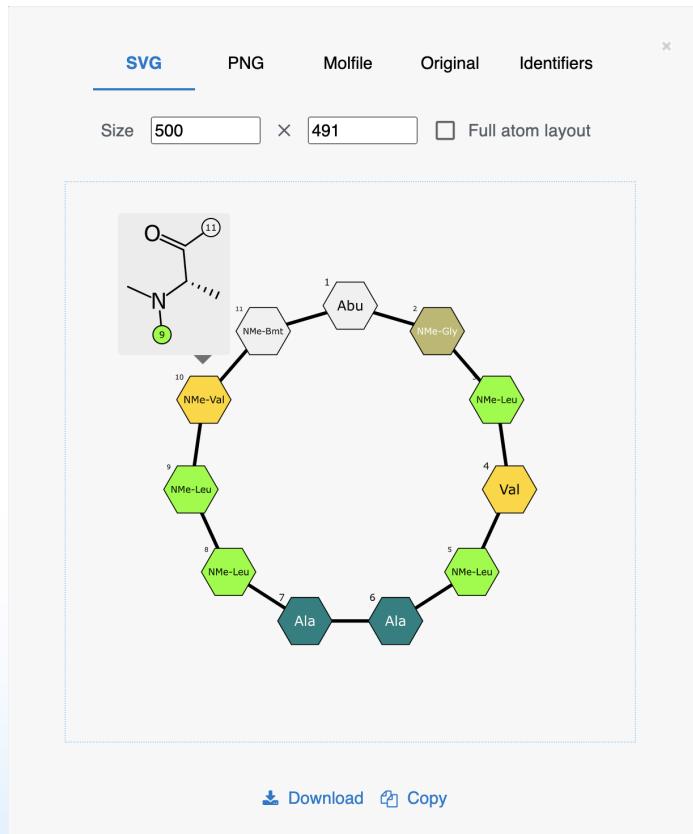
4 custom peptides,
0 custom nucleotides

A	B	
ID	Sequence	
1	19aa	GGGWGLALFKALKLPLRT[OMe] 
2	Precursor	GGGLPRT[OMe] 
3	Dimer	GGGLPRTGGGLI 
4	Trimer	GGGLPRTGGGLIATCGGLPRT 

Cyclic peptides



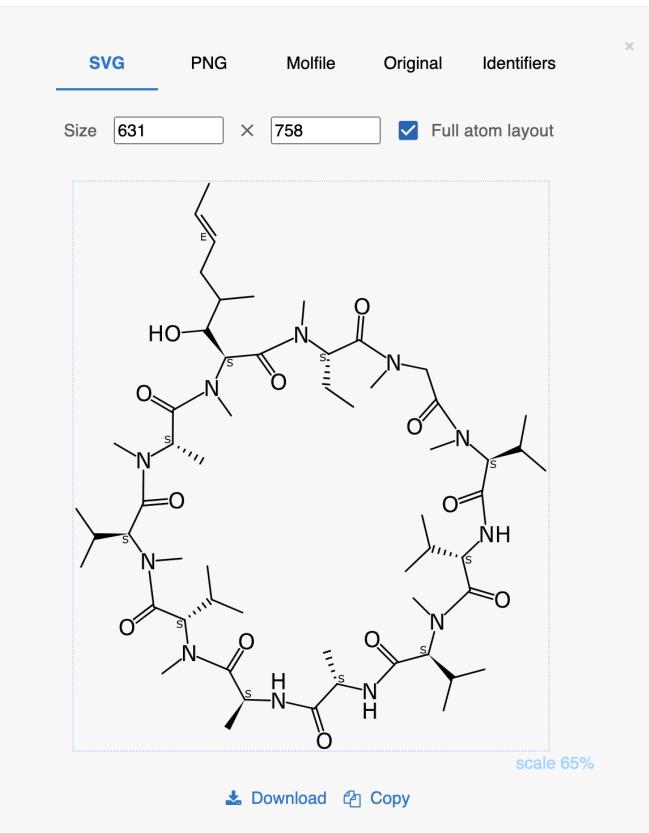
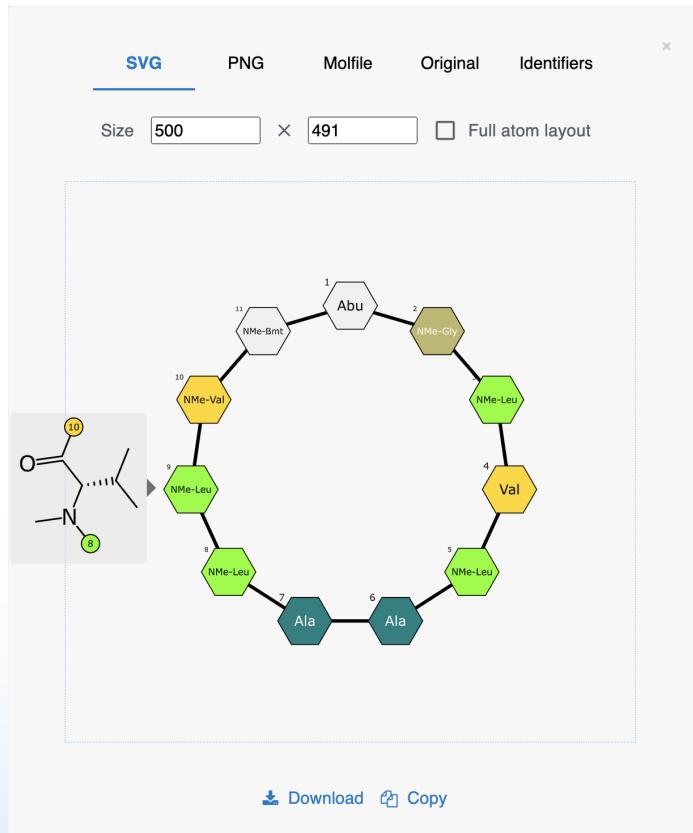
Cyclic peptides



Cyclic peptides



CDD VAULT
Complexity Simplified



Nucleotides



Explore Data ELN

Step 1: Choose Data File

File: importcodons.csv

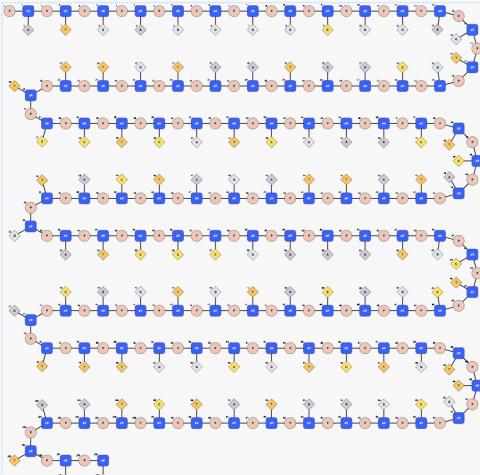
Not part of macromolecule
 Linear peptide sequence
 Cyclic peptide sequence
 Nucleotide sequence

A	Name	Sequence
1	Seq1	ATGAGGGGGAGTC TCATGATCCGAAATG ATATCTAATAG
2	Seq2	ATGAGGGGGAGTC TCATGATCCGAAATG ATAAAATAA
3	Seq3	ATGCCTTGGCCCTG CAGGAGGCAAATGCG GTGA

Preview Macromolecule

SVG PNG Molfile Original Identifiers

Size 1537 x 1541 Full atom layout



scale 32%

[Download](#) [Copy](#)

Alex Clark

Data

molecules · Owner: Alex Clark

0 custom peptides,
0 custom nucleotides

[Close](#)

Explore Data ELN Inventory Import Data Reports Settings [✉](#) Alex Clark ▾

Step 1: Choose Data File and Parse Step 2: Map Fields Step 3: Commit Data

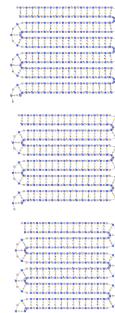
File: importcodons.csv Project: Macromolecules · Owner: Alex Clark

Compose macromolecules from columns [Edit](#)

Not part of macromolecule Backbone: RNA DNA
 Linear peptide sequence Strands: Single Double
 Cyclic peptide sequence
 Nucleotide sequence

0 custom peptides,
0 custom nucleotides

A	B
1 Seq1	Name: Sequence ATGAGGGGGAGTGCATAATGTTCTCGTGAACTCATATTAGA TCATGATCCGAAATGCTGACATGTGACATTGGCGTCTGGCGA ATATCAATAG
2 Seq2	ATGAGGGGGAGTGCATAATGTTCTCGTGAACTCATATTAGA TCATGATCCGAAATGCTGACATGTGACATTGGCGTCTGGCGA ATAAAATAA
3 Seq3	ATGCCTTGGCCCTCTAACGCCAGCAAAAGCCTTGCGCTGCA CAGGAGGCAAATGCGGAAGAATTATGCGAACATTGGATGACAC GTGA



Nucleotides



Explore Data ELN

Step 1: Choose Data File

File: importcodons.csv

Not part of macromolecule
 Linear peptide sequence
 Cyclic peptide sequence
 Nucleotide sequence

A	Name	Sequence
1	Seq1	ATGAGGGGGAGTC TCATGATCCGAAATG ATATCTAATAG
2	Seq2	ATGAGGGGGCGGAGTC TCATGATCCGAAATG ATAAAATAA
3	Seq3	ATGCCTTGGCCCTG CAGGAGGCAAATGCG GTGA

Preview Macromolecule

SVG PNG Molfile Original Identifiers

Size × Full atom layout

scale 25%

[Download](#) [Copy](#)

[Close](#)

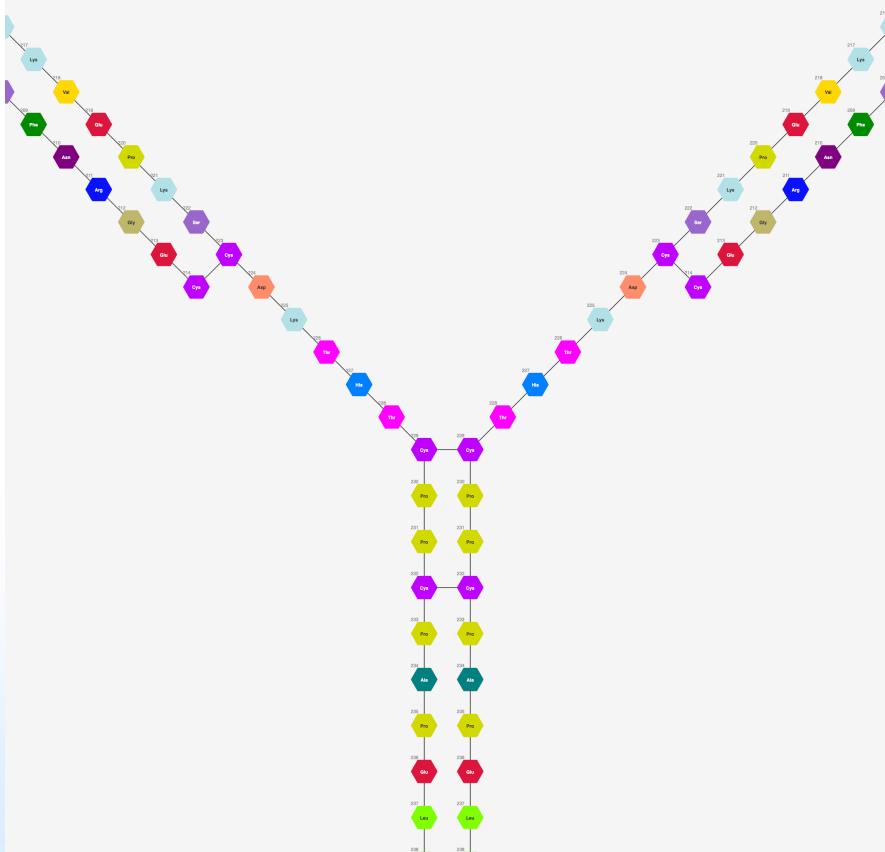
Alex Clark ▾

Data

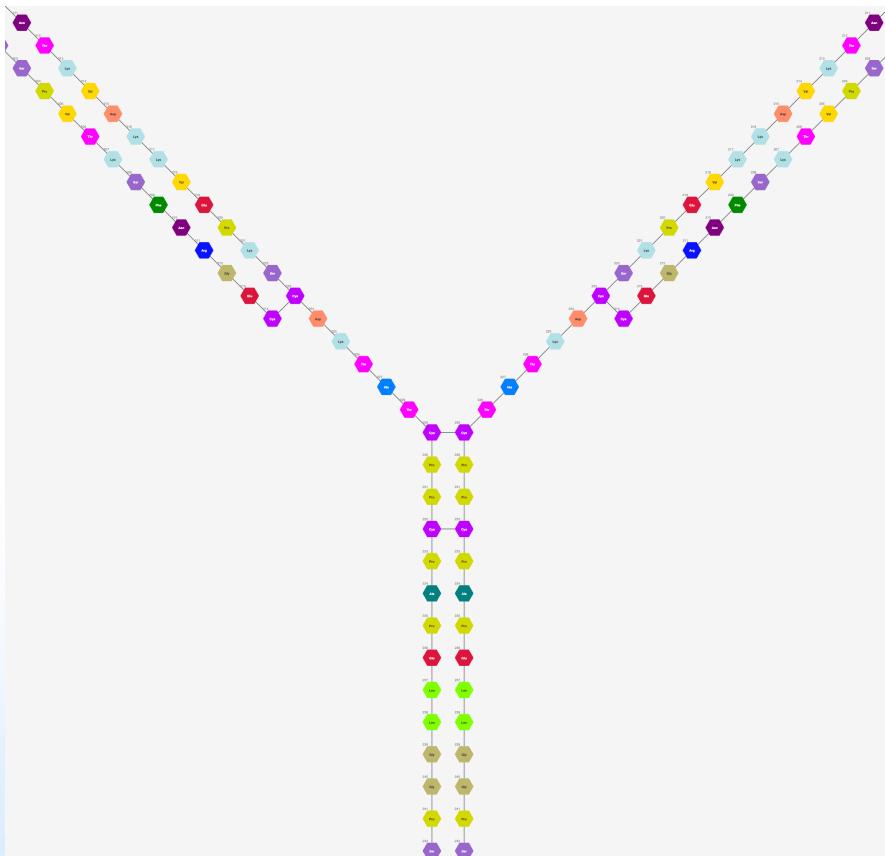
molecules · Owner: Alex Clark

0 custom peptides,
0 custom nucleotides

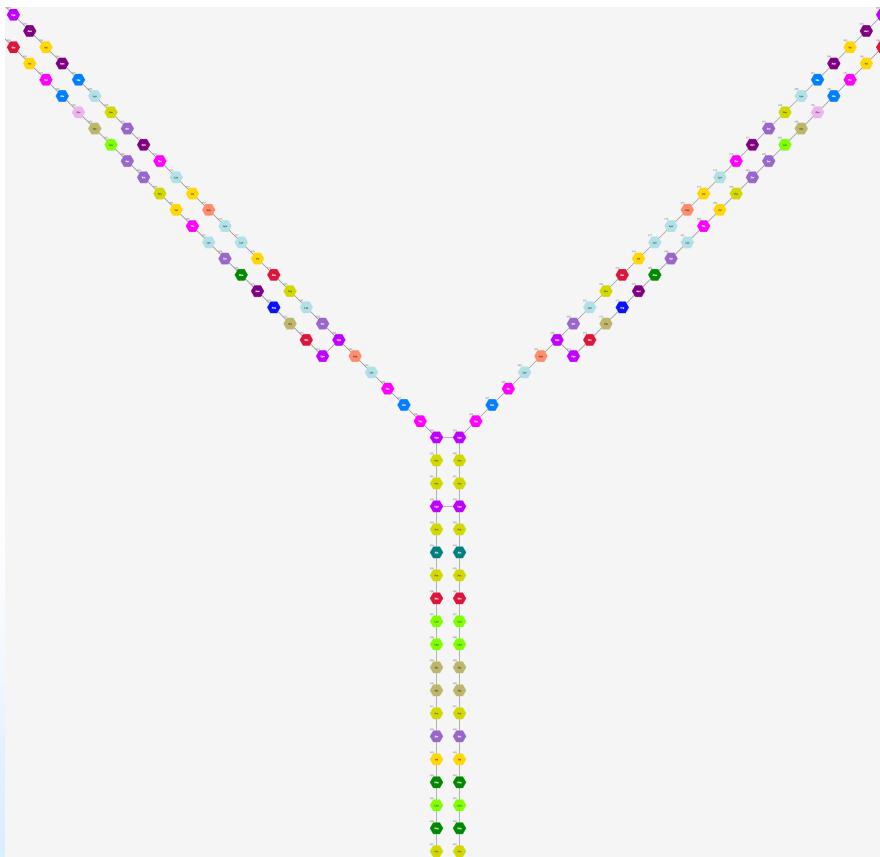
The problem of scale



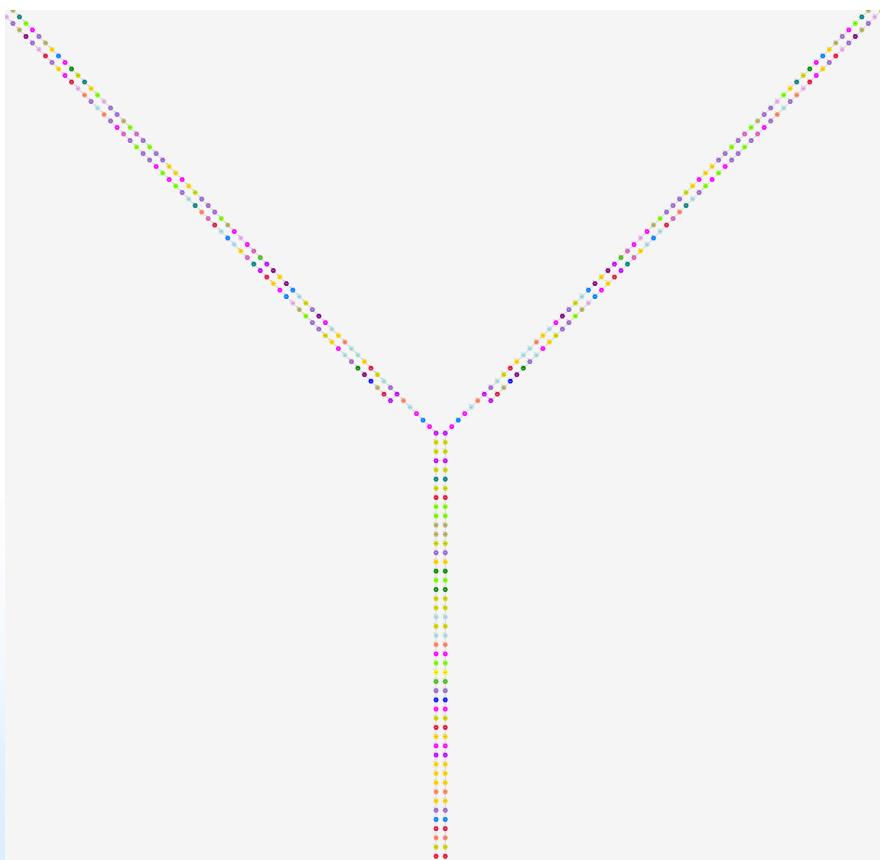
The problem of scale



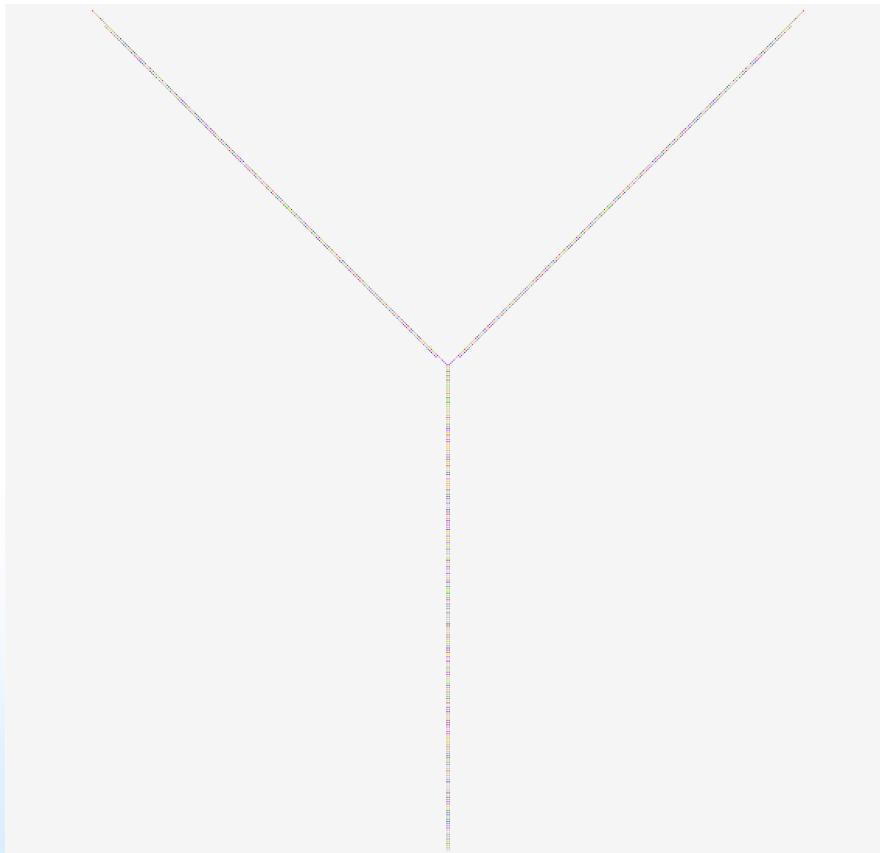
The problem of scale



The problem of scale



The problem of scale



Antibodies: Basic



Just have sequences
and drug conjugate
structures?

Not a problem:
collection of objects

Antibodies: Metadata

Light

DIQMTQSPSSLSASVGDRVTITCRASQDVNTAVAWYQQKPGKAPKLLIYSASFLYSGVPSRFSGSR
SGTDFTLTISSLQPEDFATYYCQQHYTPPTFGQGTKVEIKRTVAAPSVFIFPPSDEQLKSGTASVV
CLLNNFYPREAKVQWKVDNALQSGNSQESVTEQDSKDSTYSLSSTTLSKADYEKHKVYACEVTH
QGLSSPVTKSFNRGEC**{VL:1-107}{CL:108-214}{@214:1}**

Heavy

EVQLVESGGGLVQPGGSLRLSCAASGFNIKDTYIHWVRQAPGKGLEWVARIYPTNGYTRYADSVK
GRFTISADTSKNTAYLQMNSLRAEDTAVYYCSRWGGDGFYAMDYWGQGTLTVSSASTKGPSVF
PLAPSSKSTSGGTAAALGCLVKDYFPEPVTWSWNSGALTSGVHTFPAVLQSSGLYSLSSVTVPS
LGTQTYICNVNKHPSNTKVDKKVEPKSCDKTHTCPPCPAPELLGGPSVFLFPPKPKDTLMISRTPE
VTCVVVDVSHEDPEVKFNWYVDGVEVHNAKTKPREEQYNSTYRVVSVLTVLHQDWLNGKEYKC
KVSNKALPAPIEKTIASKAKGQPREPQVYTLPPSREEMTKNQVSLTCLVKGFYPSDIAVEWESNGQP
ENNYKTTPPVLDSDGSFFLYSKLTVDKSRWQQGNVFSCSVMHEALHNHYTQKSLSLSPGK
{VH:1-120}{CH1:121-218}{CH2:234-343}{CH3:344-450}{H:219-233}{@223:1}{@229:2}
{@232:3}{\\$264:1}{\\$324:1}{\\$370:1}{\\$428:1}



Domain Region



Disulfide Bridge



Hinge Region

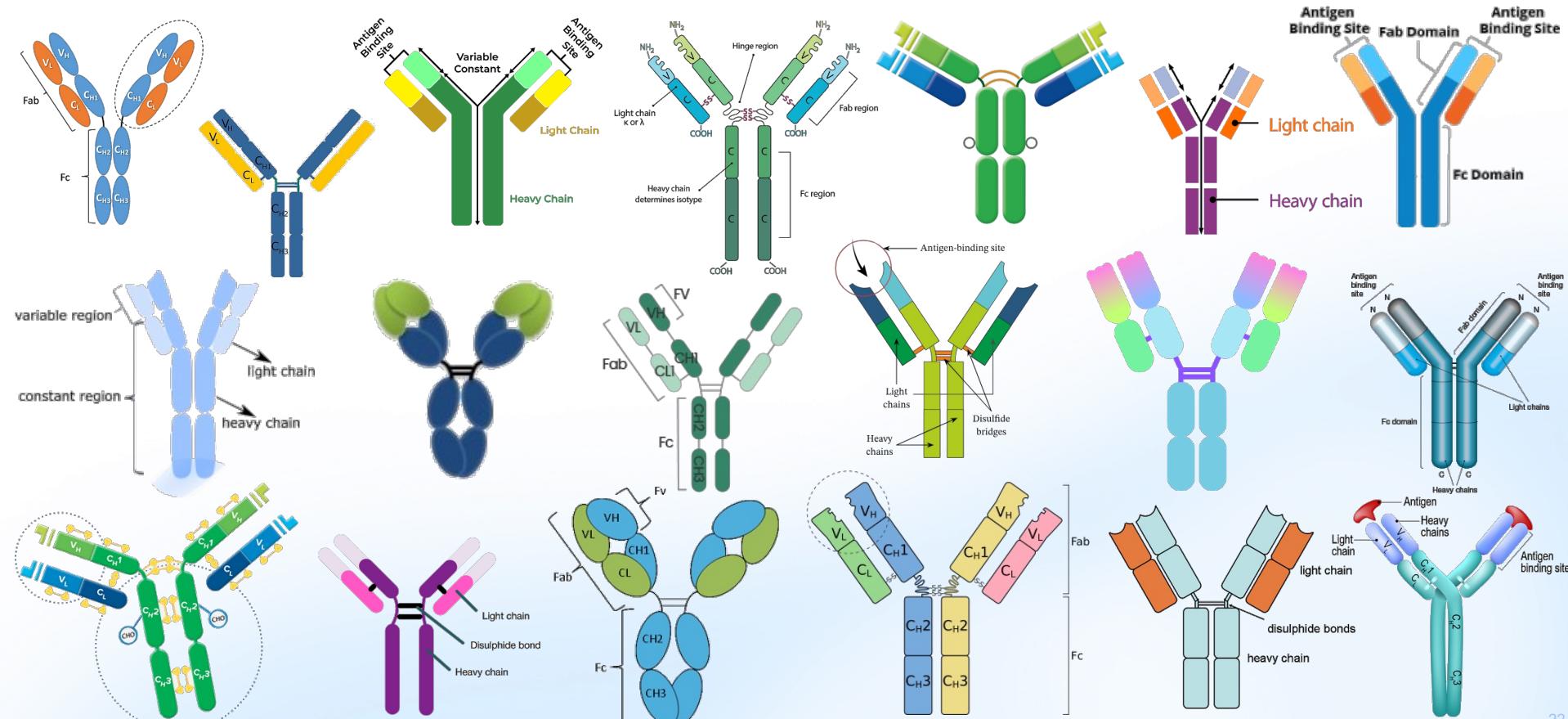


Drug Conjugate Site

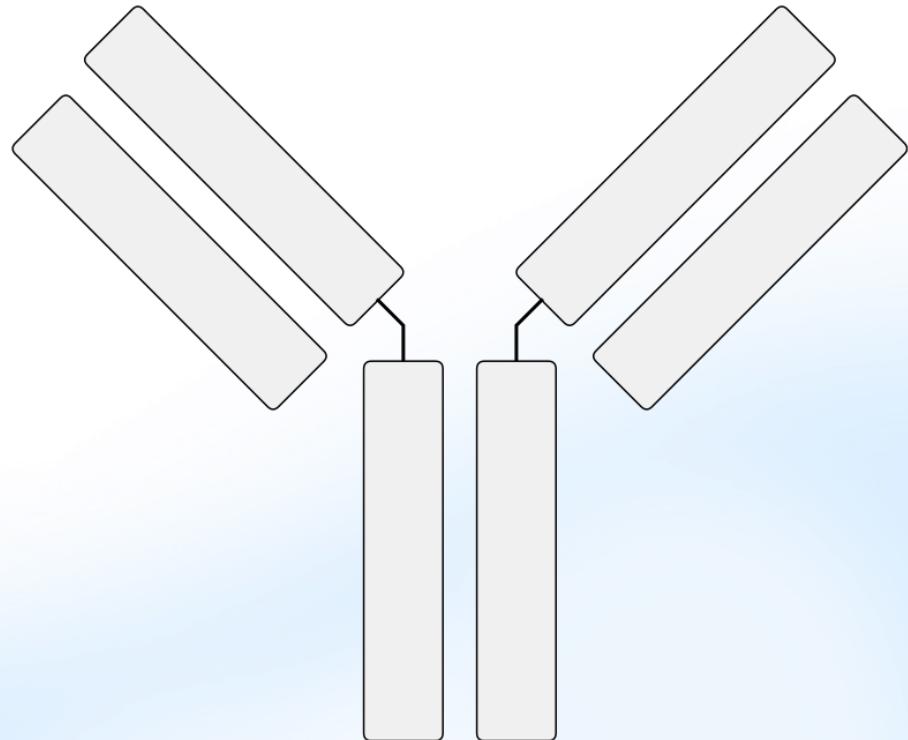
Diversity of displays



CDD VAULT
Complexity Simplified

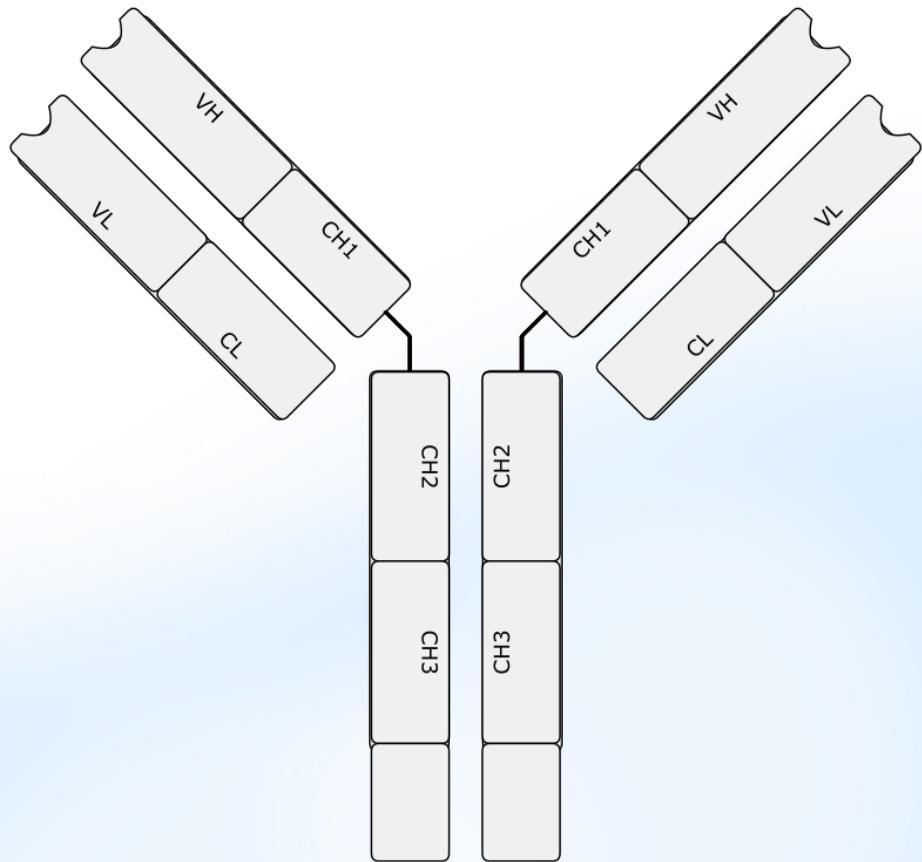


Draw basic Y-shape outline



Draw basic Y-shape outline

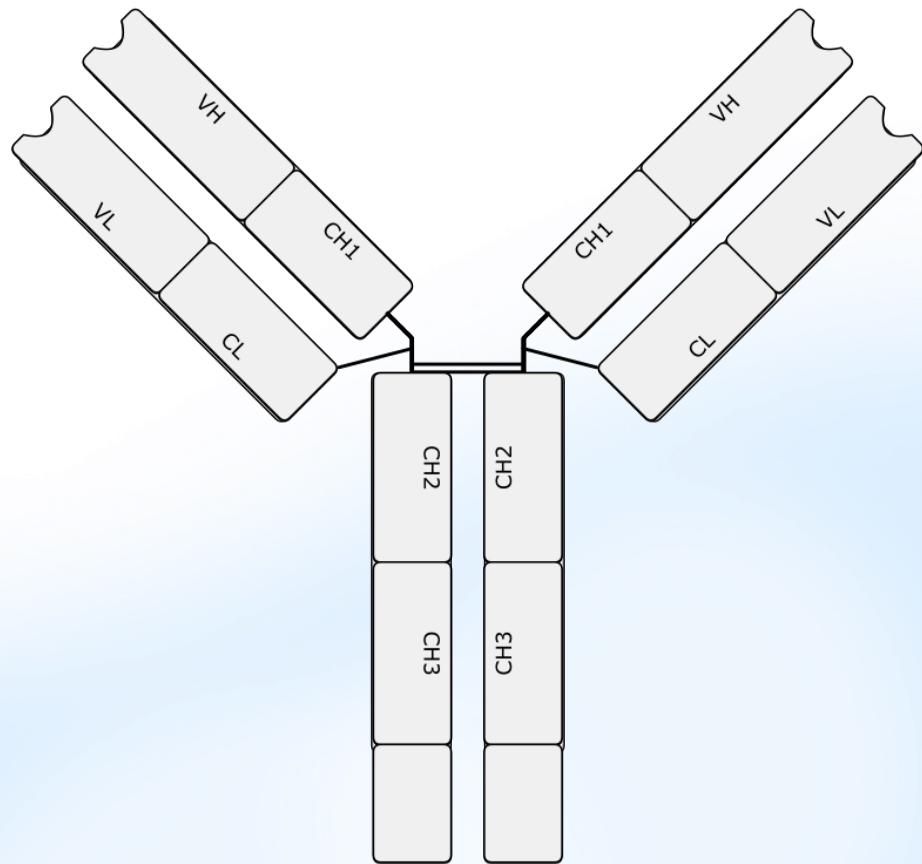
Segment domain regions



Draw basic Y-shape outline

Segment domain regions

Inter-chain disulfide links

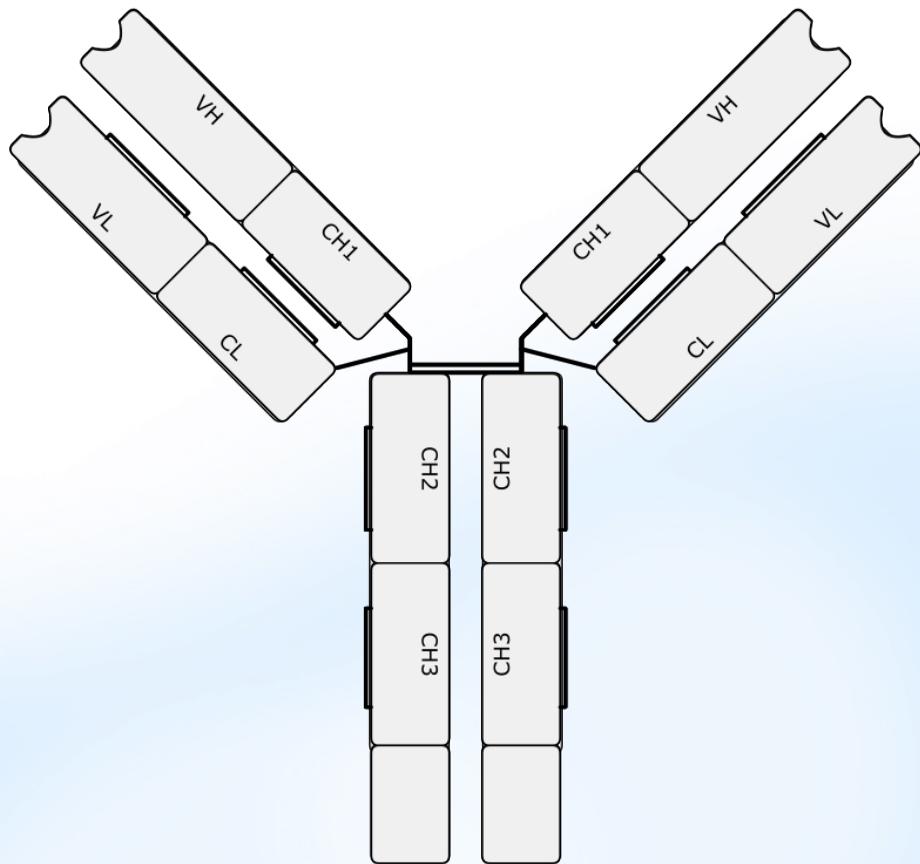


Draw basic Y-shape outline

Segment domain regions

Inter-chain disulfide links

Intra-chain disulfide links



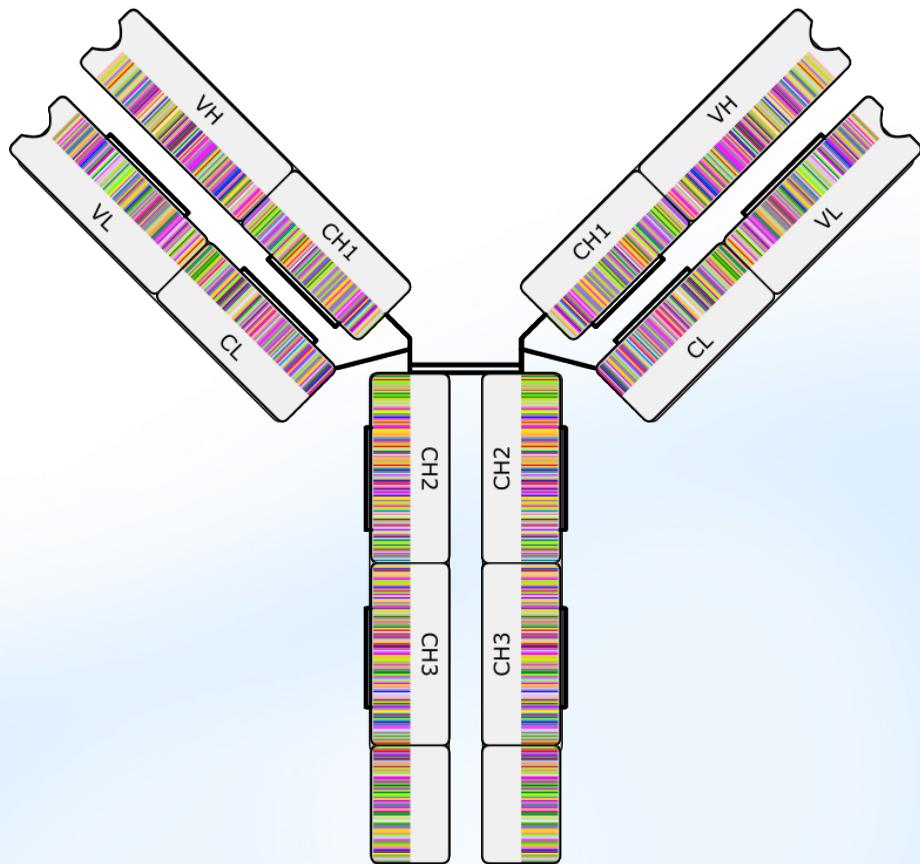
Draw basic Y-shape outline

Segment domain regions

Inter-chain disulfide links

Intra-chain disulfide links

Sequence barcoding

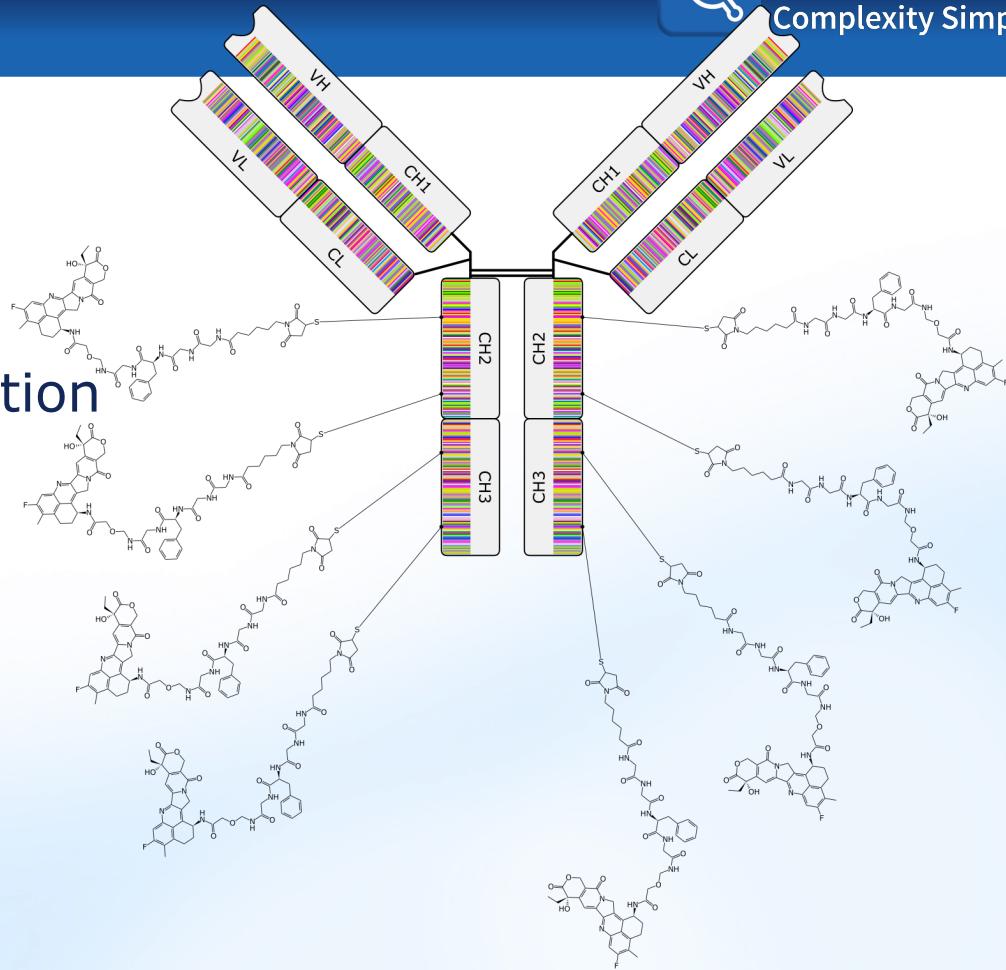


Drug conjugates

Tethered initial placements

Rigid body position optimization

Anchor point is indicated



One dedicated employee for new features/fixes: **Tad Hurst**

SCSR, MRV, stereochemistry, atropisomers...



PostgreSQL cartridge

Python microservice for i/o & layered **fingerprint**

lambdas for **standardization**

WASM for client-side

Diagrams for biologists

Sequences for bioinformaticians

Atoms & bonds for chemists

Standard formats and **open-source** tools where possible

Provide **all** information: get all the functionality

Excessive flexibility without constraints creates
unusable standards