



254th ACS National Meeting Washington Aug 2017

PubChem as a Biologics Database

Noel O'Boyle and Roger Sayle

NextMove Software

Evan Bolton

PubChem, NCBI-NIH



PUBCHEM INTERFACE



Degarelix



Download



Share



Help



+ Contents



1 2D Structure

2 3D Status

● 3 Biologic Description

4 Names and Identifiers

5 Chemical and Physical Properties

6 Related Records

7 Chemical Vendors

8 Drug and Medication Information

9 Pharmacology and Biochemistry

10 Use and Manufacturing

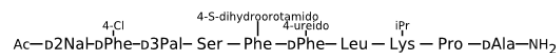
11 Safety and Hazards

12 Toxicity

3 Biologic Description



3.1 Biologic Depiction



► from PubChem

3.2 Biologic Line Notation



IUPAC Condensed

Ac-D-2NaI-D-Phe(4-Cl)-D-3Pal-Ser-Phe(4-S-dihydroorotamido)-D-Phe(4-ureido)-Leu-Lys(iPr)-Pro-D-Ala-NH₂

► from PubChem

Sequence

XXXSXXLXPA

► from PubChem

https://pubchem.ncbi.nlm.nih.gov/classification/#hid=72

PubChem Classification Browser

Browse PubChem data using a classification of interest, or search for PubChem records associated with a specific classification (e.g. Gene Ontology: DNA repair). [More...](#)

Select classification

PubChem: PubChem Compound TOC

Search selected classification

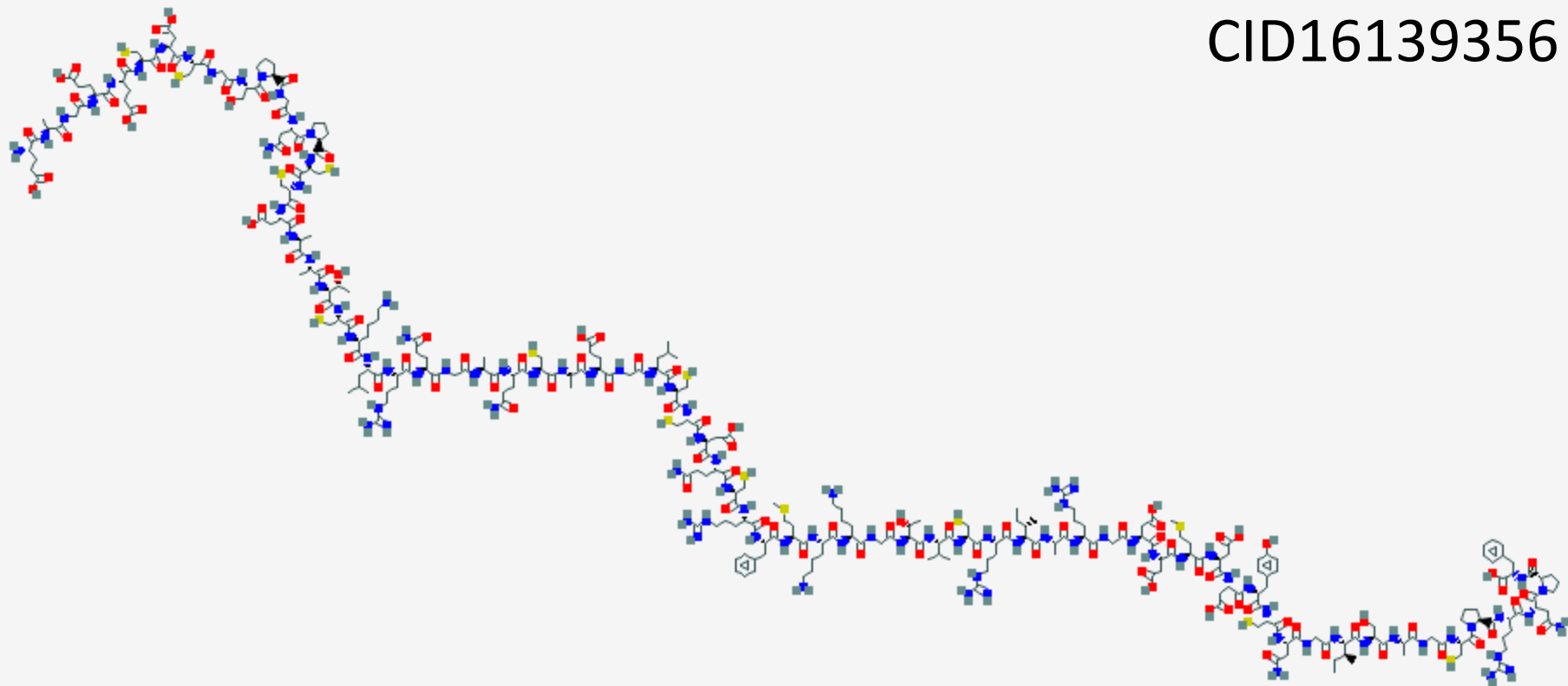
Keyword

Browse PubChem: PubChem Compound TOC Tree

▼ PubChem Compound TOC	?	29,294,437
▶ Agrochemical Information	?	1,934
▼ Biologic Description	?	475,304
Biologic Depiction	?	455,265
Biologic Line Notation	?	473,412

Click and sort by MW

CID16139356



CID16139356

H—Glu — Ala — Gly — Glu — Glu — Cys — Asp — Cys — Gly — Ser

Pro — Gly — Asn — Pro — Cys — Cys — Asp — Ala — Ala — Thr

Cys — Lys — Leu — Arg — Gln — Gly — Ala — Gln — Cys — Ala

Glu — Gly — Leu — Cys — Cys — Asp — Gln — Cys — Arg — Phe

Met — Lys — Lys — Gly — Thr — Val — Cys — Arg — Ile — Ala

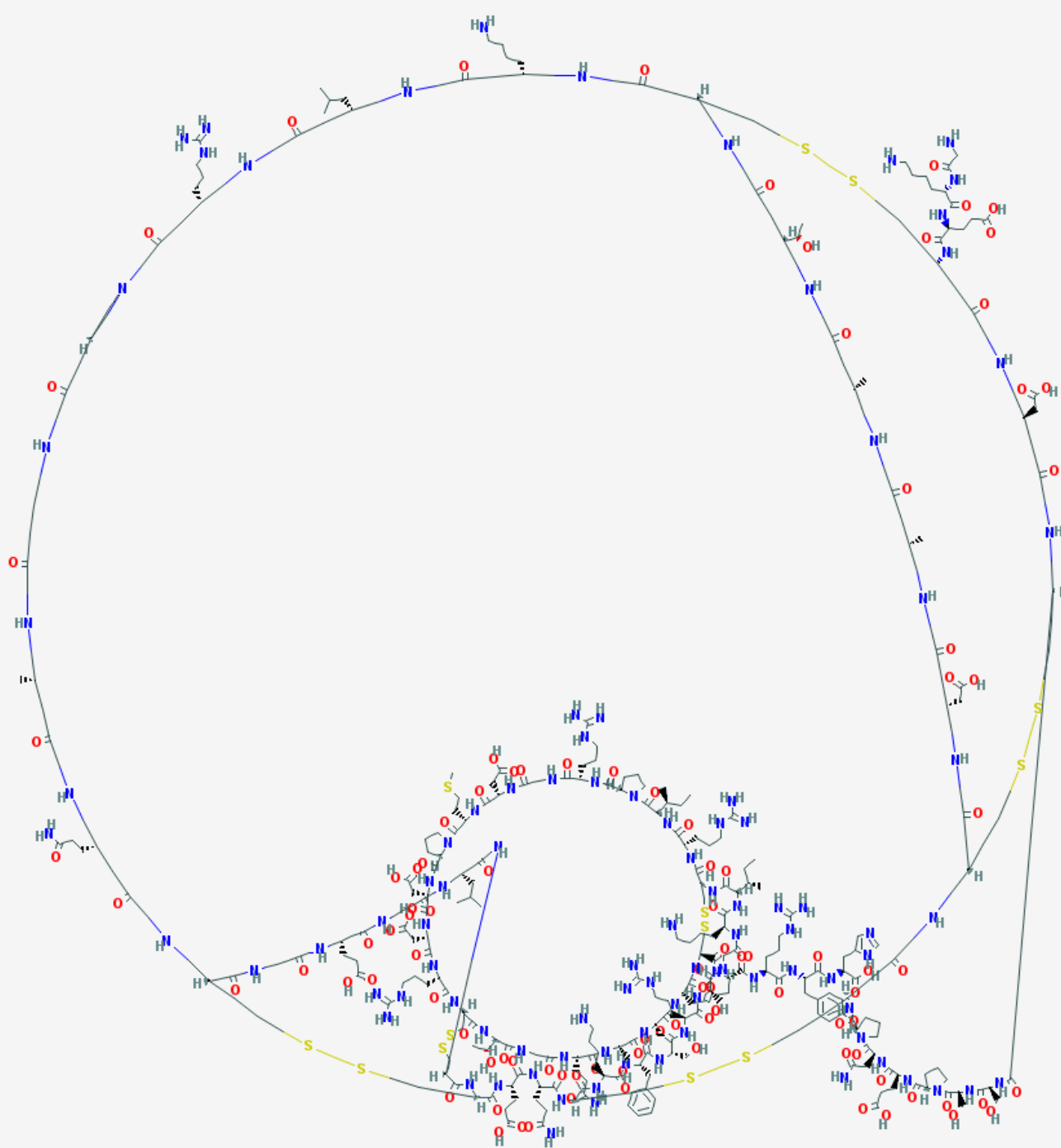
Arg — Gly — Asp — Asp — Met — Asp — Asp — Tyr — Cys — Asn

Gly — Ile — Ser — Ala — Gly — Cys — Pro — Arg — Asn — Pro

Phe—OH



CID56842075 Rhodostomin



H—Gly — Lys — Glu — Cys — Asp — Cys — Ser — Ser — Pro — Glu

CID56842075
Rhodostomin

Asn — Pro — Cys — Cys — Asp — Ala — Ala — Thr — Cys — Lys

Leu — Arg — Pro — Gly — Ala — Gln — Cys — Gly — Glu — Gly

Leu — Cys — Cys — Glu — Gln — Cys — Lys — Phe — Ser — Arg

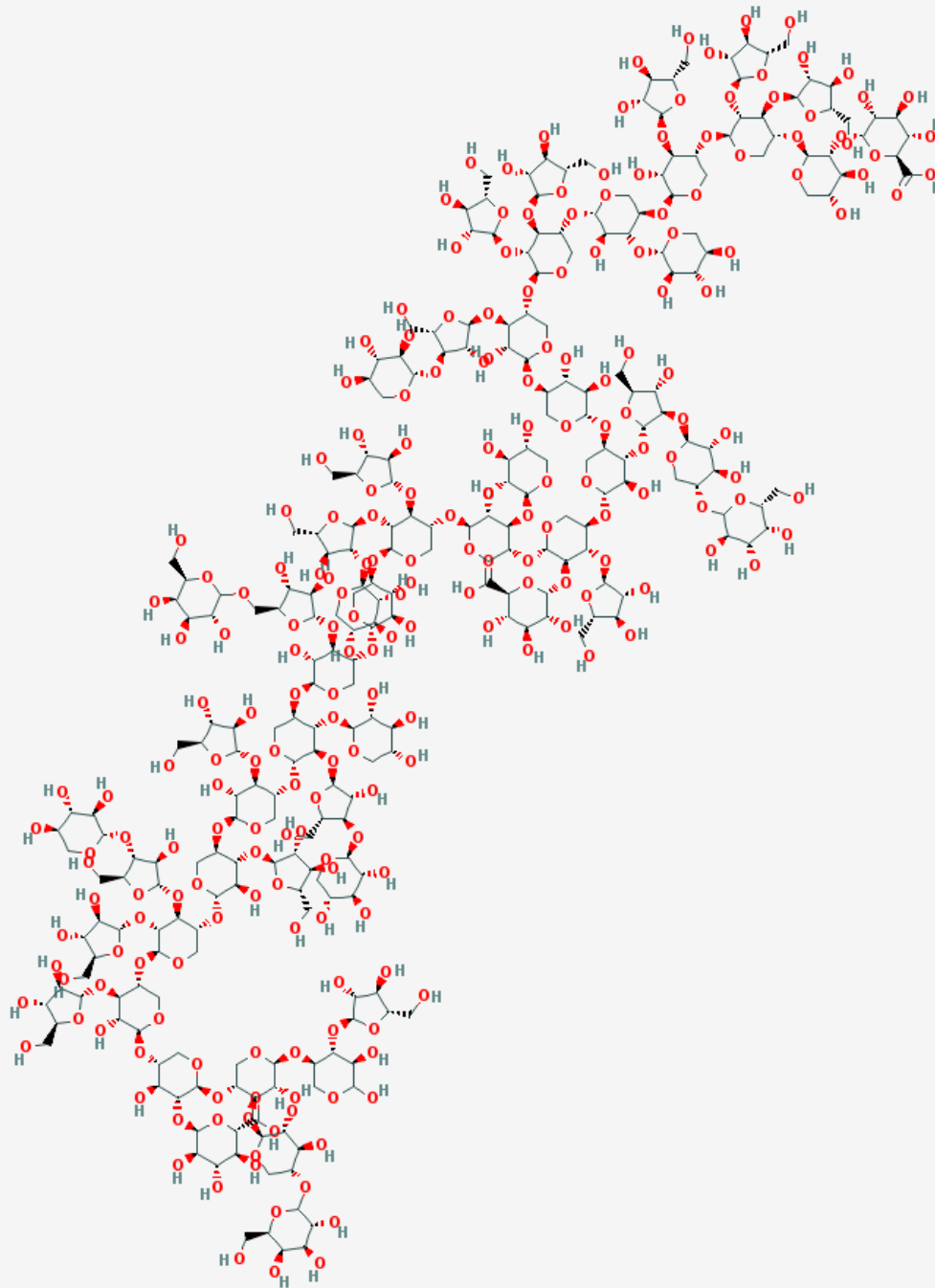
Ala — Gly — Lys — Ile — Cys — Arg — Ile — Pro — Arg — Gly

Asp — Met — Pro — Asp — Asp — Arg — Cys — Thr — Gly — Gln

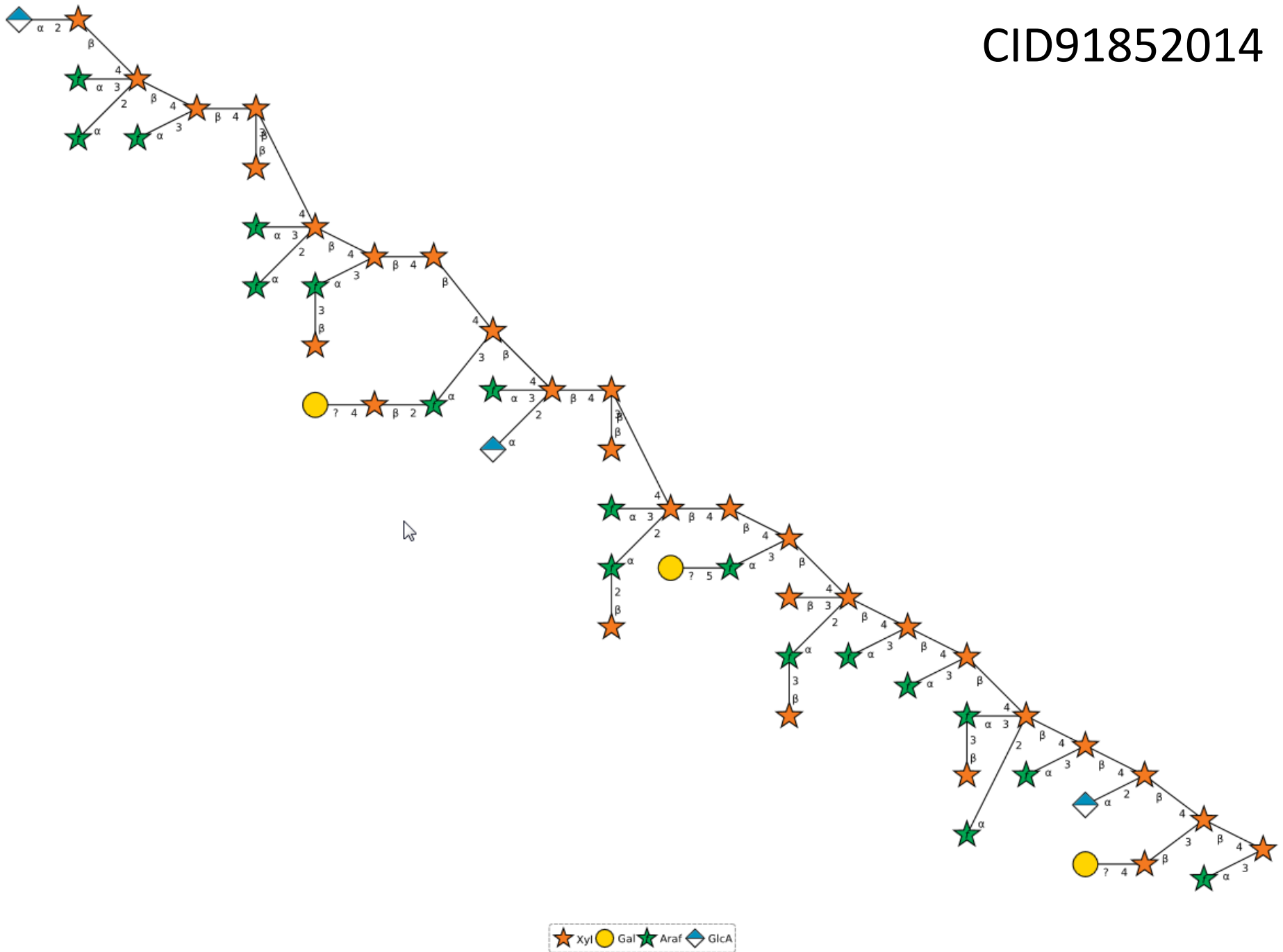
Ser — Ala — Asp — Cys — Pro — Arg — Tyr — His—OH



CID91852014



CID91852014

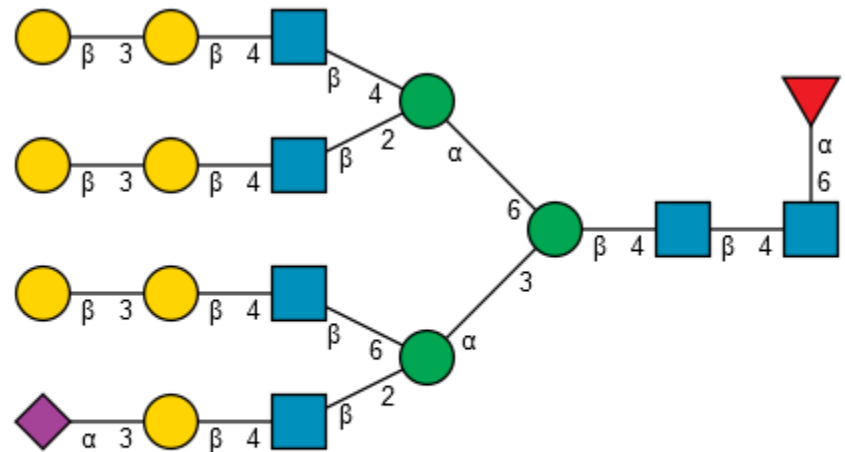
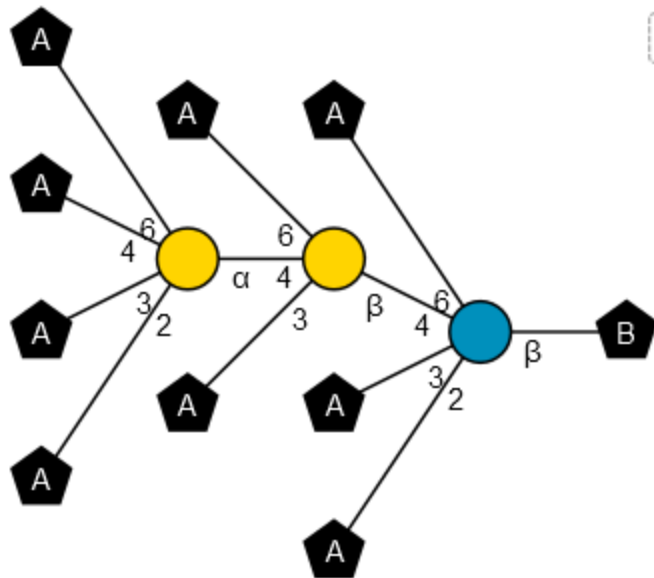
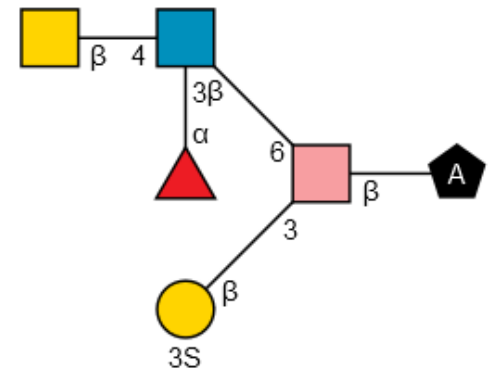
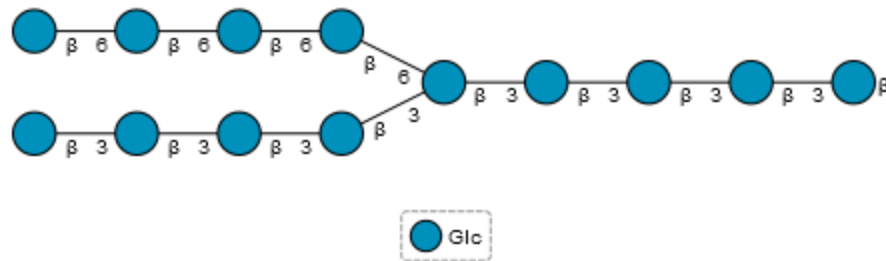
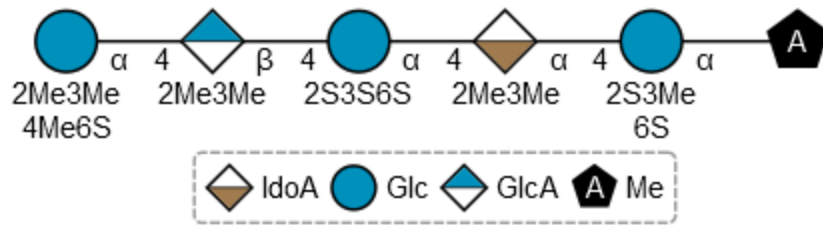


PUBCHEM – A SMALL MOLECULE DATABASE?

- People don't think of PubChem as a **peptide** database
 - ~110K X-rays of proteins in PDB
 - ~500K peptides in PubChem
- People don't think of PubChem as a **saccharide** database
 - ~80K oligosaccharides in GlyTouCan
 - ~67K oligosaccharides in PubChem



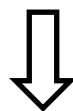
531,618 contain a monosaccharide, of which 66,740 can be depicted



HOW MANY MONOSACCHARIDES PRESENT?

113 aldoses, ketoses, aldonic and uronic acids with from 5-9 carbons

AltA, Glc, L-Man, L-Gal, Fru, L-gro-D-glcHept



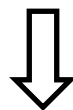
407 including deoxy variants, ring variants

L-Glcf, Mans, 2-deoxy-D-manHept, 3-deoxy-D-glcOct2ulo-onic



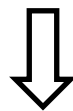
971 including anomeric stereo

a-Man, 3,4-deoxy-a-D-eryHex, b-Tyv



7094 including common substituents at non-anomeric positions

Xylf5Me, a-L-ManNAc3Ac4Ac6Ac, Glc2P3P6P

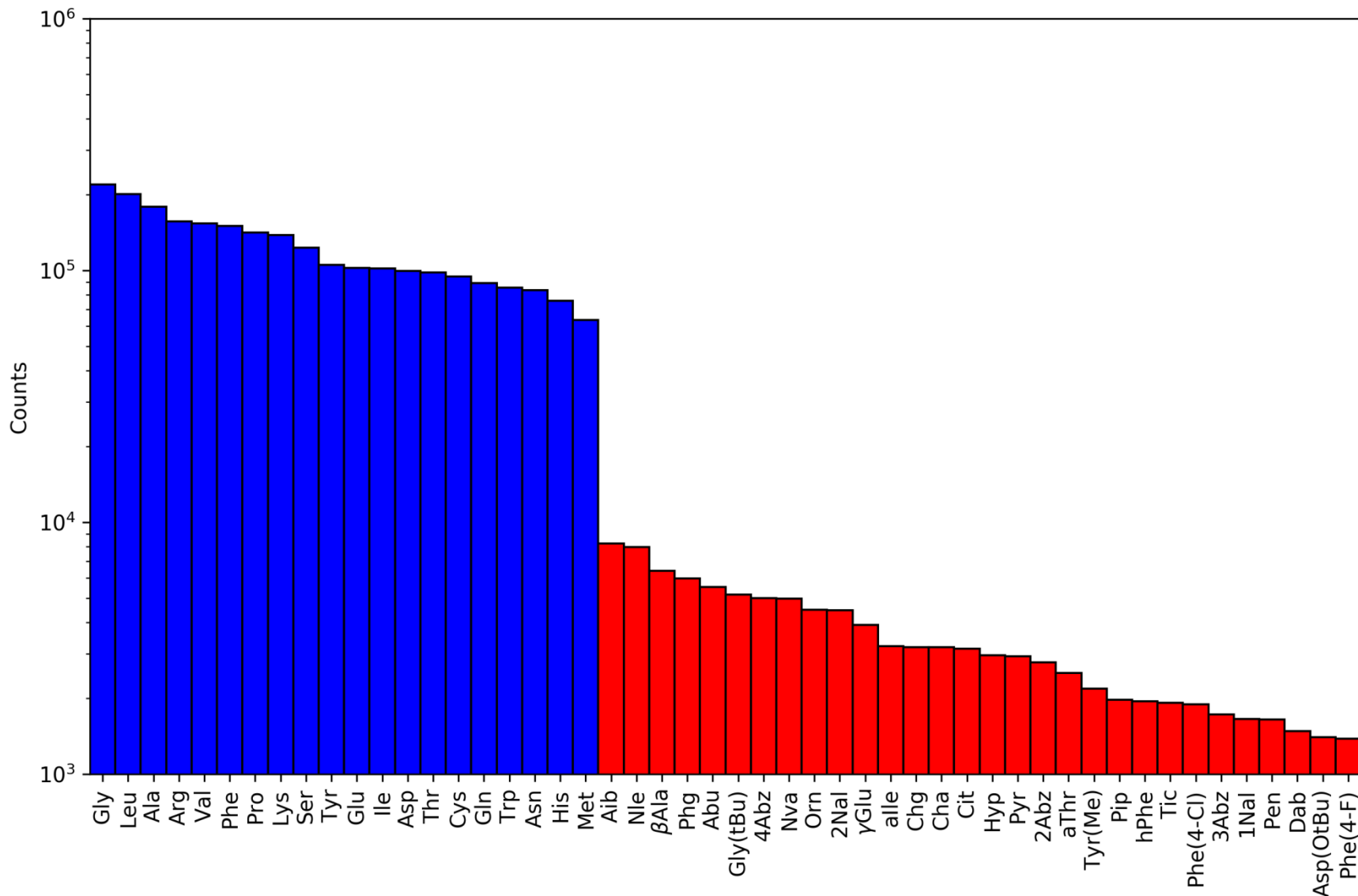


26641 including any substituent anywhere

Bz(-2)[Tos(-3)]Ara4Ac(b)-O-Me, TMS(-4)[TMS(-6)]GlcNAc3Me(a)-O-Me



AMINO ACIDS IN PUBCHEM STRUCTURES CONTAINING AT LEAST THREE AMINO ACIDS



HOW MANY AMINO ACIDS PRESENT?

20 common amino acids

Ala, Cys, Lys, Thr



87 amino acids

Ala, Cys, Hcy, Lys, 2Nal, Ncy, Thr



1095 including substituents

Thr, Thr(*t*Bu), Thr(Bn), Thr(PO₃H₂)



3546 including stereo variants, terminal variants, linker variants, α -methylated

Thr, D-Thr, DL-Thr, aThr, Thr-ol, aMeThr



8125 including N-substituted variants

Thr, Me-Thr, Boc-Thr, Me₂-Thr, Fmoc-N(Me)Thr



HOW MANY PEPTIDES PRESENT?

- Depends how you count...
- 447,026 have 3 or more amino acids
- 668,229 structures are recognised in their entirety as peptides



<div> <div>iturelix</div> <div>CID16130938</div> </div>	<div> <div> <div>4-Cl</div> <div>nicotinoyl</div> <div>nicotinoyl</div> <div>iPr</div> </div> <div> <div>Ac</div> <div>d2Nal</div> <div>dPhe</div> <div>d3Pal</div> <div>Ser</div> <div>Lys</div> <div>dLys</div> <div>Leu</div> <div>Lys</div> <div>Pro</div> <div>dAla</div> <div>NH₂</div> </div> </div>
<div> <div>elamipretide</div> <div>CID11764719</div> </div>	<div> <div> <div>2,6-diMe</div> </div> <div> <div>H</div> <div>dArg</div> <div>Tyr</div> <div>Lys</div> <div>Phe</div> <div>NH₂</div> </div> </div>
<div> <div>histrelin</div> <div>CID25077993</div> </div>	<div> <div> <div>1-Bn</div> </div> <div> <div>H</div> <div>Pyr</div> <div>His</div> <div>Trp</div> <div>Ser</div> <div>Tyr</div> <div>dHis</div> <div>Leu</div> <div>Arg</div> <div>Pro</div> <div>NHEt</div> </div> </div>
<div> <div>icatibant</div> <div>CID71364</div> </div>	<div> <div> <div>H</div> <div>dArg</div> <div>Arg</div> <div>Pro</div> <div>Hyp</div> <div>Gly</div> <div>2Thi</div> <div>Ser</div> <div>dTic</div> <div>Oic</div> <div>Arg</div> <div>OH</div> </div> </div>
<div> <div>linacлотide</div> <div>CID16158208</div> </div>	<div> <div> <div>H</div> <div>Cys</div> <div>Cys</div> <div>Glu</div> <div>Tyr</div> <div>Cys</div> <div>Cys</div> <div>Asn</div> <div>Pro</div> <div>Ala</div> <div>Cys</div> </div> <div> <div>Thr</div> <div>Gly</div> <div>Cys</div> <div>Tyr</div> <div>OH</div> </div> </div>
<div> <div>valinomycin</div> <div>CID5649</div> </div>	<div> <div> <div>dAla</div> <div>dVal</div> <div>Val</div> <div>dVal</div> <div>dAla</div> <div>dVal</div> <div>Val</div> <div>Val</div> <div>Ala</div> <div>Val</div> </div> <div> <div>Val</div> <div>Val</div> </div> </div>

SEQUENCE REPRESENTATION

- Depending on the task or quality of datasource, **different sequence representations** may be preferred
 - Distinguish between D-/L-/DL- amino acids using upper/lowercase?
 - D-Ala as a or A
 - Distinguish between sidechain stereo variants?
 - alloThr as X or T
 - Distinguish between substituted amino acids and their parent?
 - Ser(PO₃H₂) as X or S



EXACT SEQUENCE SEARCH

- Given that features of the structure are normalized or ignored
 - Exact sequence search can be used to find similar structures (sequence as hash)
- Create **hierarchy of similarity**
 - First, those structures with the same sequence, if we normalise as much as possible
 - Then successively discriminate based on stereochemistry, side-chain substitution



SEARCH FOR KEMPTIDE: LRRASLG

H-Leu — Arg — Arg — Ala — Ser — Leu — Gly—OH

LRRASLG

LRRASLG

Ac-DL-Leu-DL-Arg-DL-Arg-DL-Ala-DL-Ser-DL-Leu-Gly-OH	78069426	acetyl-kemptide
Ac-Leu-Arg-Arg-Ala-Ser-Leu-Gly-OH	71429096	acetyl-kemptide
H-DL-Leu-DL-Arg-DL-Arg-DL-Ala-DL-Ser-DL-Leu-Gly-NH ₂	85062657	kemptide amide
H-DL-Leu-DL-Arg-DL-Arg-DL-Ala-DL-Ser-DL-Leu-Gly-OH	100074	kemptide
H-DL-Leu-DL-Arg-DL-Arg-DL-Ala-DL-Ser-DL-Leu-Gly-OH.TFA	118797564	
H-Leu-Arg-Arg-Ala-Ser-Leu-Gly-NH ₂	9897033	kemptide amide
H-Leu-Arg-Arg-Ala-Ser-Leu-Gly-OH	9962276	kemptide
Unk-Leu-Arg-Arg-Ala-Ser-Leu-Gly-OH	11650926,101224399,101878757	

LRRAXLG

H-Leu-Arg-Arg-Ala-Ser(PO ₃ H ₂)-Leu-Gly-NH ₂	102212089	[Ser(PO ₃ H ₂)-5]kemptide amide
H-Leu-Arg-Arg-Ala-Ser(PO ₃ H ₂)-Leu-Gly-OH	13783725	[Ser(PO ₃ H ₂)-5]kemptide

LRraSLG

LRraSLG

H-Leu-Arg-D-Arg-D-Ala-Ser-Leu-Gly-OH	53393688	[D-Arg ₃ ,D-Ala ₄]kemptide
--------------------------------------	----------	---

LrRASLG

LrRASLG

H-Leu-D-Arg-Arg-Ala-Ser-Leu-Gly-OH	99864041	[D-Arg ₂]kemptide
------------------------------------	----------	-------------------------------

lRRASLG

lRRASLG

H-D-Leu-Arg-Arg-Ala-Ser-Leu-Gly-OH	99864040	[D-Leu ₁]kemptide
------------------------------------	----------	-------------------------------

lrRASLG

lrRASLG

H-D-Leu-D-Arg-Arg-Ala-Ser-Leu-Gly-OH	99864042	[D-Leu ₁ ,D-Arg ₂]kemptide
--------------------------------------	----------	---

DISULFIDE BRIDGING PATTERNS

- Use a sequence representation to find peptides with different **disulfide bridges**
 - Does not occur naturally
 - Either errors by depositor, or artificially created
- Convert peptides with at least four cysteines to sequence format and collate
 - 16 cases found with different bridges
 - 12 were erroneous, 4 real



DISULFIDE BRIDGING PATTERNS

I**C**CNPAC**C**GPKYSC

CID11480353



CID101041637



CHEMICAL BIOLOGY & DRUG DESIGN



[Explore this journal >](#)

Controlled syntheses of natural and disulfide-mispaired regioisomers of α -conotoxin SI[†]

B. Hargittai, G. Barany

First published: December 1999 [Full publication history](#)

DOI: 10.1034/j.1399-3011.1999.00127.x [View/save citation](#)

Deposited by Nikajii



DISULFIDE BRIDGING PATTERNS

GLPRKIL**CA**IAKKKGK**CK**GPLKLV**CKC**

CID71597277

H—Gly — Leu — Pro — Arg — Lys — Ile — Leu — Cys — Ala — Ile — Ala — Lys — Lys — Lys — Gly — Lys — Cys — Lys — Gly — Pro

Leu — Lys — Leu — Val — Cys — Lys — Cys—OH

CID71597445

H—Gly — Leu — Pro — Arg — Lys — Ile — Leu — Cys — Ala — Ile — Ala — Lys — Lys — Lys — Gly — Lys — Cys — Lys — Gly — Pro

Leu — Lys — Leu — Val — Cys — Lys — Cys—OH

CID71597707

H—Gly — Leu — Pro — Arg — Lys — Ile — Leu — Cys — Ala — Ile — Ala — Lys — Lys — Lys — Gly — Lys — Cys — Lys — Gly — Pro

Leu — Lys — Leu — Val — Cys — Lys — Cys—OH

Deposited by NIAID (National Institute of Allergies and Infectious Diseases)

DISULFIDE BRIDGING PATTERNS

GLPRKIL**CA**IAKKKGK**CK**GPLKLV**CK**C

CID71597277

H—Gly — Leu — Pro — Arg — Lys — Ile — Leu — Cys — Ala — Ile — Ala — Lys — Lys — Lys — Gly — Lys — Cys — Lys — Gly — Pro —

— Leu — Lys — Leu — Val — Cys — Lys — Cys—OH

CID71597445

H—Gly — Leu — Pro — Arg — Lys — Ile — Leu — Cys — Ala — Ile — Ala — Lys — Lys — Lys — Gly — Lys — Cys — Lys — Gly — Pro —

Amino Acids (2012) 43:751–761

DOI 10.1007/s00726-011-1125-6

ORIGINAL ARTICLE

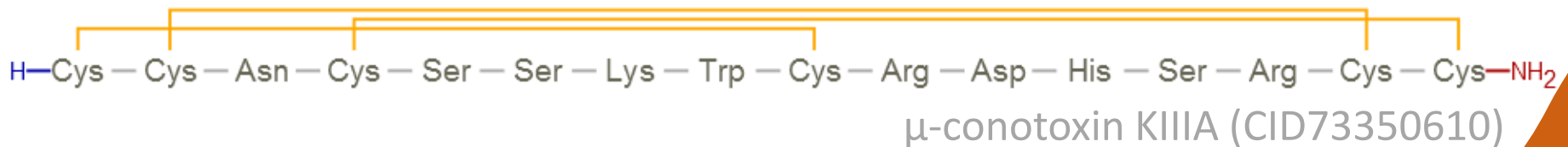
Lasiocepsin, a novel cyclic antimicrobial peptide from the venom of eusocial bee *Lasioglossum laticeps* (Hymenoptera: Halictidae)

Lenka Monincová • Jiřina Slaninová • Vladimír Fučík •
Oldřich Hovorka • Zdeněk Voburka • Lucie Bednářová •
Petr Maloň • Jitka Štokrová • Václav Čerovský

Deposited by NIAID (National Institute of Allergies and Infectious Diseases)

KNOTTINS

- Peptides with three disulfide bridges, where one threads through the macrocycle formed by the others (see KNOTTIN Database <http://knottin.cbs.cnrs.fr>)
- Interesting leads for drug discovery
 - Stable fold, sequence tolerant, small



- A necessary (but not sufficient) condition is an arrangement of Cys bridges 123123
 - 90 examples in PubChem



SEQUENCE VARIATION

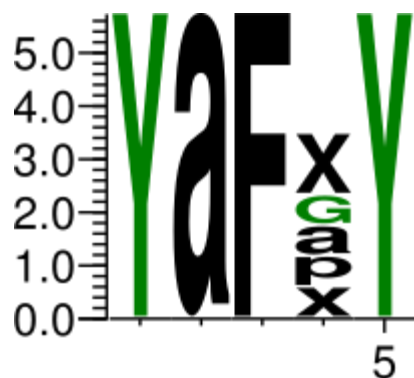
- Which parts of a sequence have seen the most **variants**?
 - Of interest for drug discovery, activity modulation
- Looked for variants of a sequence that are one substitution away from a known peptide
 - Required strict matching of the sequence to minimize 'mutations' due to errors



PubChem

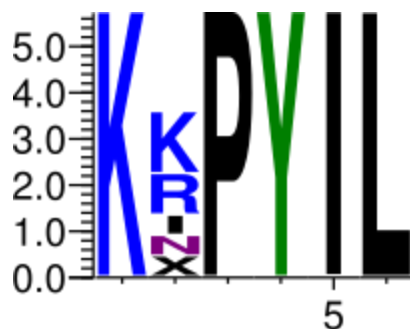
Casokefamide

YaFaY



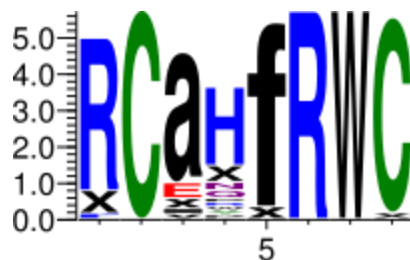
Neuromedin N

KIPYIL



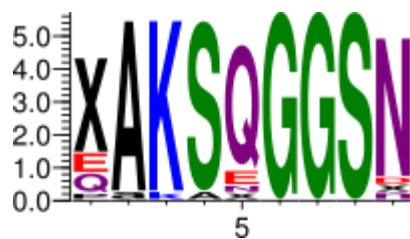
Setmelanotide

RCaHfRWC



Thymulin

XAKSQGGSN

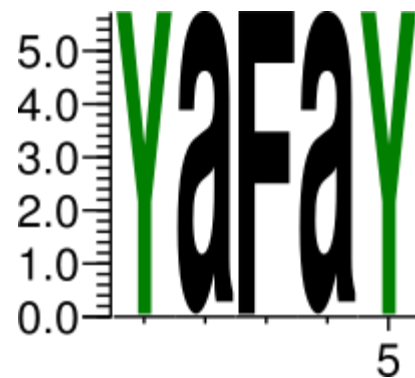
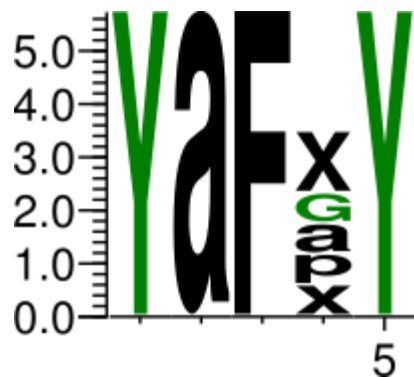


PubChem

ChEMBL

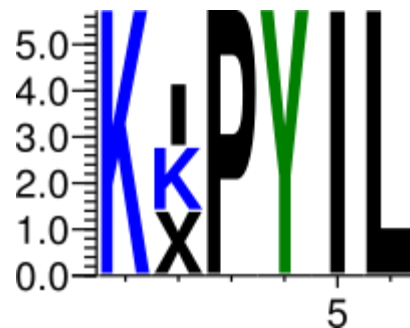
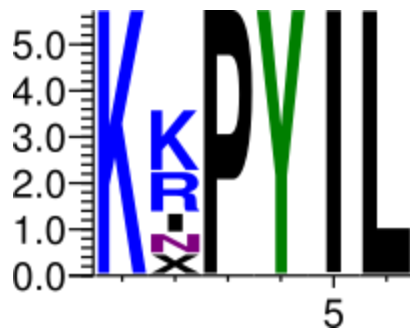
Casokefamide

YaFaY



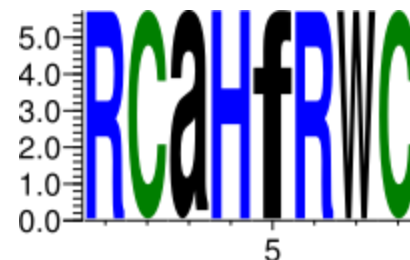
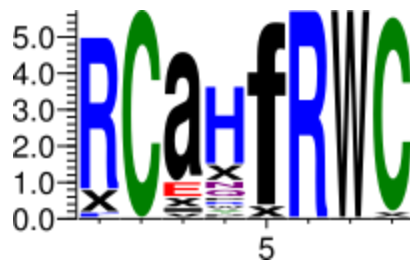
Neuromedin N

KIPYIL



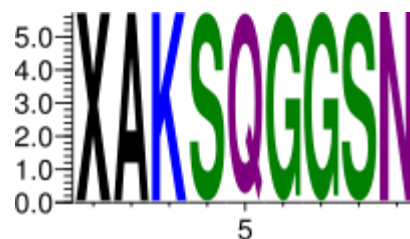
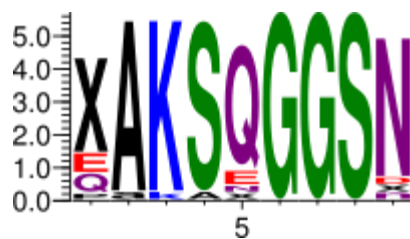
Setmelanotide

RCaHfRWC



Thymulin

XAKSQGGSN



ARE PUBCHEM PEPTIDES VARIANTS OF KNOWN PEPTIDES?

- Hypothesis: observed peptides are close variants of a small number of known peptides
- Curated database of oligopeptides of biological interest (currently 452 entries)
- 10.5% of the 170,708 peptides of length 5 or greater in PubChem can be named as variants of these
 - argipressin (1-8)
 - Cbz-cholecystokinin octapeptide (2-7) amide
 - [Ile1,Ser2,Ser8]cyphokinin



SUMMARY

- PubChem is a rich source of information on **oligopeptides** and **oligosaccharides**
 - Often heavily modified, rather than natural
- Due to chemical modifications, we need to think in terms of **10s of thousands** of monomers
- Sequence representations act as a key:
 - To collate similar peptides
 - To find sequence variants, sites of variation
 - To find disulfide bridge variants, knottins



ACKNOWLEDGEMENTS

- Paul Thiessen PubChem

SUGAR & SPLICE



<http://nextmovesoftware.com>

<http://nextmovesoftware.com/blog>

noel@nextmovesoftware.com





Rank Sales 2016	Trade Name	Name	Type of biologic
1	Humira	adalimumab	Monoclonal antibody
2	Harvoni	ledipasvir/sofosbuvir	
3	Enbrel	etanercept	Protein attached to monoclonal antibody
4	Rituxan	rituximab	Monoclonal antibody
5	Remicade	infliximab	Monoclonal antibody
6	Revlimid	lenalidomide	
7	Avastin	bevacizumab	Monoclonal antibody
8	Herceptin	trastuzumab	Monoclonal antibody
9	Lantus	insulin glargine	Protein
10	Pprevnar 13	Pneumococcal vaccine	Polysaccharides attached to carrier protein
11	Xarelto	rivaroxaban	
12	Eylea	aflibercept	Protein attached to monoclonal antibody
13	Lyrica	pregabalin	
14	Neulasta	pegfilgrastim	PEG attached to protein
15	Advair Diskus	fluticasone/salmeterol	

Source: Genetic Engineering & Biotechnology News

<http://www.genengnews.com/the-lists/the-top-15-best-selling-drugs-of-2016/77900868>



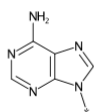
REPRESENTATIONS AT VARIOUS LEVELS

- All-atom
- Hydrogen-suppressed graphs
 - How to handle non-standard valencies, cycles
- Consider monomer as superatom
 - How to handle cycles, monomer variants such as modifications and opposite stereo

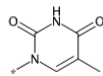


MOST COMMON UNRECOGNISED SUBSTITUENTS ON SUGARS

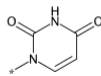
- Increased number of perceived oligosaccharides from 51,273 to 66,740



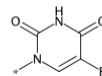
21434



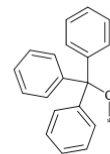
12372



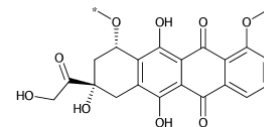
9853



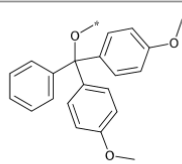
2091



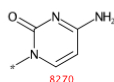
1963



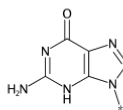
1760



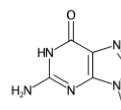
9438



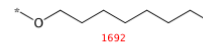
8270



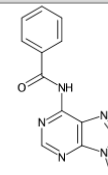
5373



1757



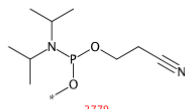
1692



1569



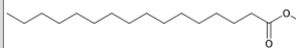
2956



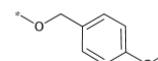
2779



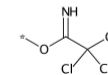
2713



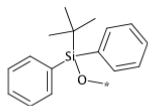
1516



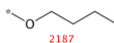
1487



1480



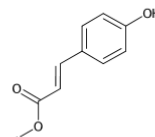
2230



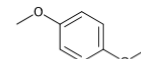
2187



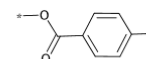
2118



1478



1473



1435