**SciFinder Discovery Platform**

# TRAINING

Karin Färber  KFaerber@acs-i.org
Ernst Aichinger
ACS International / CAS
March 2022

CAS
A division of the
American Chemical Society

# Content curation & Digitization

**Delivering decades of paper-to-digital scientific content creation expertise**



**CAS SciFinder$^n$**



Source Selection → Text Indexing → Reaction Indexing → Substance Indexing → Data Model & Structure → Search Architecture → App Development

# Components and Search Types

## SciFinder Discovery Platform



**CAS SciFinder$^n$ search options**

- References
- Visualizations — Chemscape / Bioscape
- Biosequences
- Suppliers
- Retrosynthesis Planner
- Reactions
- Properties and formulas
- Substances

### Detailed Formulations

## CAS Formulus

Quick access to the detailed formulation information for all active ingredients and excipients.
Main areas: pharma, agrochemistry and cosmetics.

### Detailed Analytical Methods

## CAS Analytical Methods

A single source for searching and comparing the latest published scientific analytic methods. Search and compare hundreds of thousands of methods across multiple scientific fields of study.

CAS
A division of the
American Chemical Society

# CAS CONTENT COLLECTION

**Curated, connected data analyzed by scientific experts fluent in over 50 languages**

| | |
|---|---|
| CAS REGISTRY | CAS Commercial Sources |
| CAS Reactions | CAS Formulations |
| CAS References | Medline (largest portion of PubMed) |
| CAS Markush | CAS Biosequences |

CAS

# Boolean operators

**CAS SciFinder$^n$ searching**

**AND** requires both concepts to be present within the document

 A AND B

**OR** requires either one or both concepts to be present

 A OR B

**NOT** excludes documents from an answer set

 A NOT B

Example: menthol **and** (food **or** candy) **not** cigarette

CAS
A division of the
American Chemical Society

# Phrase search and wildcards

Terms enclosed in double quotes will be searched in the input order

- E.g.: "transcription factor"

Wildcards for internal and right-hand truncation

- * replaces any number of characters | plasticiz* → plasticizer, plasticized, plasticizing
- ? replaces 0 or 1 character | 1,?-Dibromobutane → 1,1-Dibromobutane ... 1,4-Dibromobutane

CAS
A division of the
American Chemical Society

SEARCHING WITH SCIFINDER$^N$

# BIOSEQUENCES

# Biosequence Search Options

- **Compounds with max. 252 non-H atoms**
  - Indexed as "usual" small molecules
  - Can be searched by sequence, structure, name, molecular formula, etc.

- **Defined compounds with more than 252 non-H atoms**
  - Are not indexed with structure: NOT searchable by structure
  - Might be indexed with molecular formula
  - Can thus be searched by name/identifier, molecular formula and sequence

- **Unspecific biomolecules**
  - Can be searched by name/identifier
  - Include reference search with broad/generic definitions

CAS
A division of the
American Chemical Society

# Biosequence Search Options

- **Compounds with max. 252 non-H atoms**
  - Indexed as "usual" small molecules
  - Can be searched by sequence, structure, name, molecular formula, etc.

  Search shorter sequences for modifications and uncommon amino acids

- **Defined compounds with more than 252 non-H atoms**
  - Are not indexed with structure: NOT searchable by structure
  - Might be indexed with molecular formula
  - Can thus be searched by name/identifier, molecular formula and sequence

  BLAST
  CDR
  Motif

- **Unspecific biomolecules**
  - Can be searched by name/identifier
  - Include reference search with broad/generic definitions

  E.g. for classes of sequences or trade-names

  Retrieve sequence-related information via text search

# Biosequence Searching
## New Data, New Options

**Comprehensive collection:**

- Over 600M patent-sequence relationships from more than 1.1 M patents and 62+ patent authorities

- Over 500M sequences submitted to the non-redundant NCBI protein and nucleotide databases

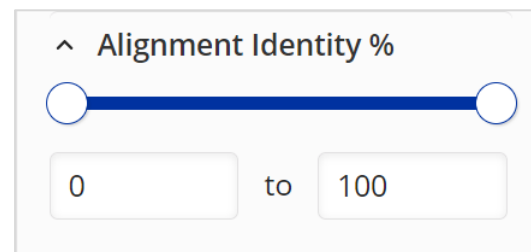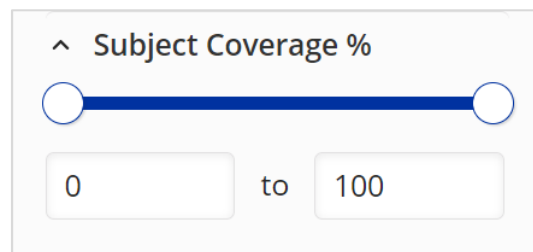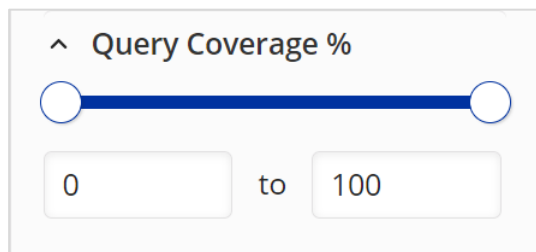- Manually curated sequences not found in electronic sequence listings and other databases

**New biosequence search option:**

- **BLAST** Search nucleotides and proteins

- **CDR** Search complementarity-determining regions of antibodies and T-Cell receptors

- **Motif** Sequence matching with wild cards and other features

**Bioscape visualization tool provides additional analysis options**

CAS
A division of the
American Chemical Society

# Alignment parameters
## Coverage and Sequence Identity percentages



| ^ Query Coverage % | ^ Subject Coverage % | ^ Alignment Identity % |
|---|---|---|
| 0 to 100 | 0 to 100 | 0 to 100 |

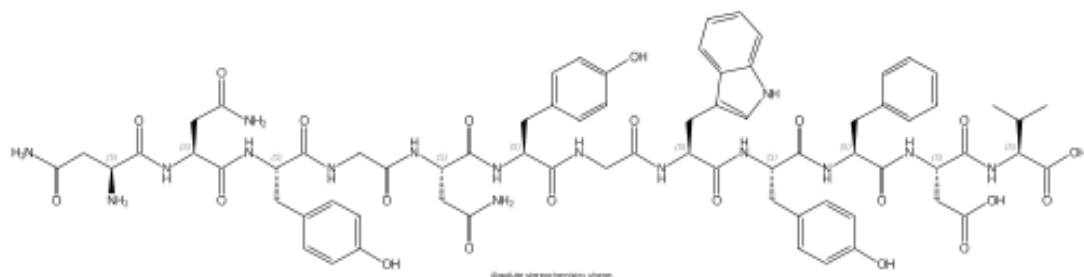Query Sequence    QQLLVVEE**G**G

Subject Sequence  QQLLVVEE**I**GS

Alignment

$$\text{Query Coverage (100\%)} = \frac{Alignment\ Length}{Query\ Length} = \frac{10}{10}$$

$$\text{Subject Coverage (91\%)} = \frac{Alignment\ Length}{Subject\ Length} = \frac{10}{11}$$

$$\text{Alignment Identity (90\%)} = \frac{Number\ of\ Matches}{Alignment\ Length} = \frac{9}{10}$$

CAS
A division of the
American Chemical Society

# Biosequences search example
**CAS RN:** 1967777-53-2

NNYGNYGWYFDV



Protein/Peptide Sequence
Sequence Length: 12

Expand All | Collapse All

▲ Other Names and Identifiers

Canonical SMILES

O=C(O)CC(NC(=O)C(NC(=O)C(NC(=O)C(NC(=O)CNC(=O)C(NC(=O)C(NC(=O)CNC(=O)C(NC(=O)C(NC(=O)C(N)CC(=O)N)CC(=O)N)CC1=CC=C(O)C=C1)CC(=O)N)CC2=CC=C(O)C=C2)CC3=CNC=4C=CC=CC43)CC5=CC=C(O)C=C5)CC=6C=CC=CC6)C(=O)NC(C(=O)O)C(C)C

Isomeric SMILES

C([C@@H](C(N[C@@H](CC1=CC=C(O)C=C1)C(N[C@@H](CC2=CC=CC=C2)C(N[C@H](C(N[C@@H](C(C)C)C(O)=O)=O)CC(O)=O)=O)=O)=O)NC(CNC([C@@H](CC3=CC=C(O)C=C3)NC([C@@H](NC(CNC([C@H](CC4=CC=C(O)C=C4)NC([C@@H](NC([C@H](CC(N)=O)N)=O)CC(N)=O)=O)=O)=O)CC(N)=O)=O)=O)=O)C=5C=6C(NC5)=CC=CC6

12 Other Names for this Substance

L-Asparaginyl-L-asparaginyl-L-tyrosylglycyl-L-asparaginyl-L-tyrosylglycyl-L-tryptophyl-L-tyrosyl-L-phenylalanyl-L-α-aspartyl-L-valine (ACI)

11: PN: WO2019053613 SEQID: 15 claimed protein

3: PN: US20160215059 SEQID: 3 claimed protein

3: PN: WO2018025221 SEQID: 3 claimed protein

3: PN: WO2019171294 SEQID: 3 claimed protein

3: PN: WO2019229614 SEQID: 3 claimed protein

3: PN: WO2020031087 SEQID: 3 claimed protein

3: PN: WO2020086476 SEQID: 3 claimed protein

3: PN: WO2020086479 SEQID: 3 claimed protein

3: PN: WO2021046289 SEQID: 3 claimed protein

3: PN: WO2021046293 SEQID: 3 claimed protein

5: PN: WO2021043961 SEQID: 3 claimed protein

CAS
A division of the
American Chemical Society

# Flexible search examples

# BLAST results
## View alignments, filters, links to patents

# BLAST results
## Export sequence data to EXCEL

# BioScape Visualization
## Interactive view of sequence answer set

**Position**
Similarity to Query
(multiple dimensions)

**Color**
Similarity measured
by Alignment Identity

**Query position**
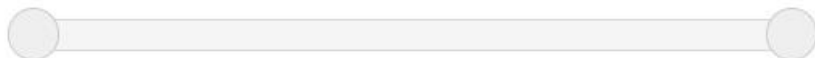Indicated by
turquoise circle ◎

# CDR Search
## Filters Options
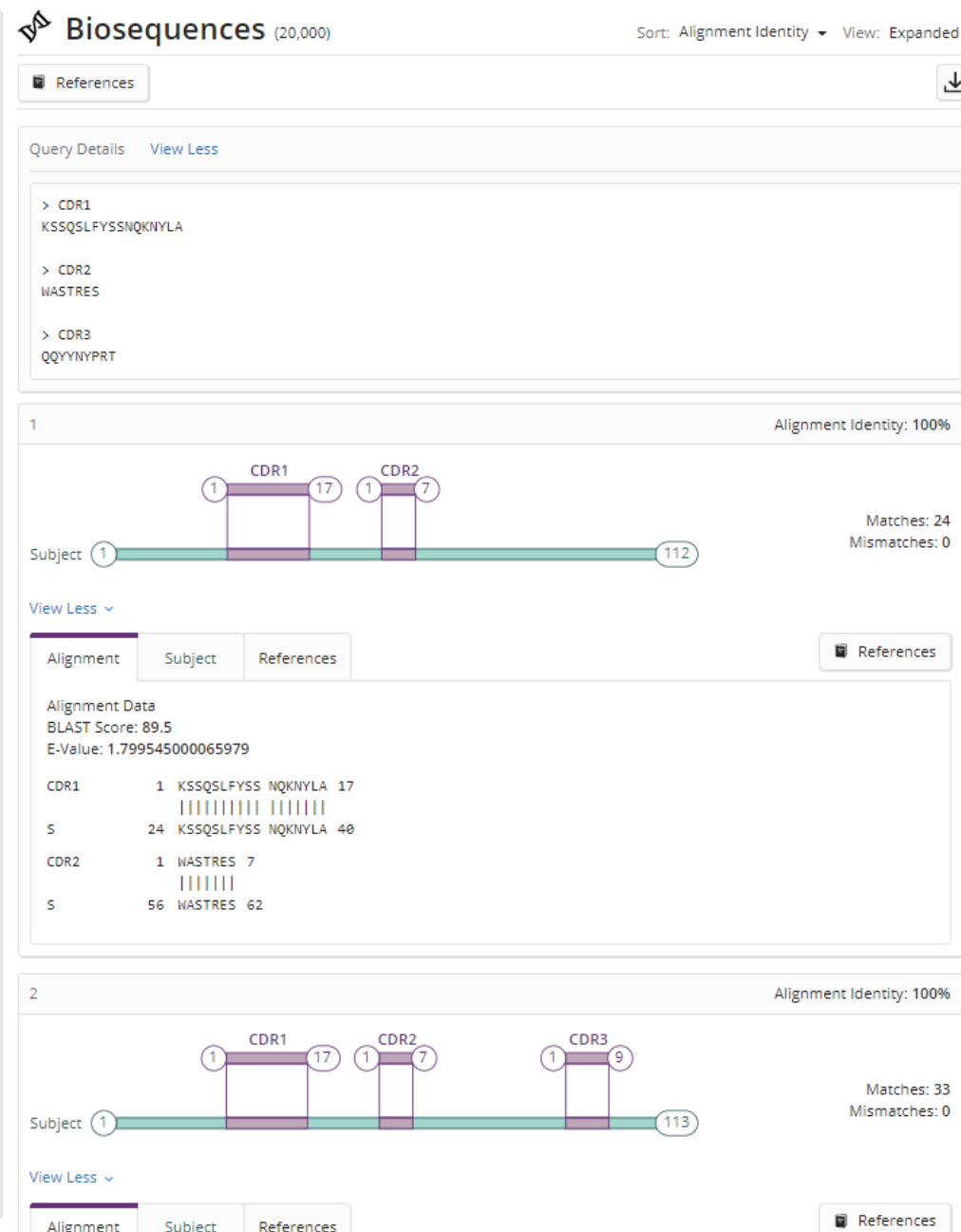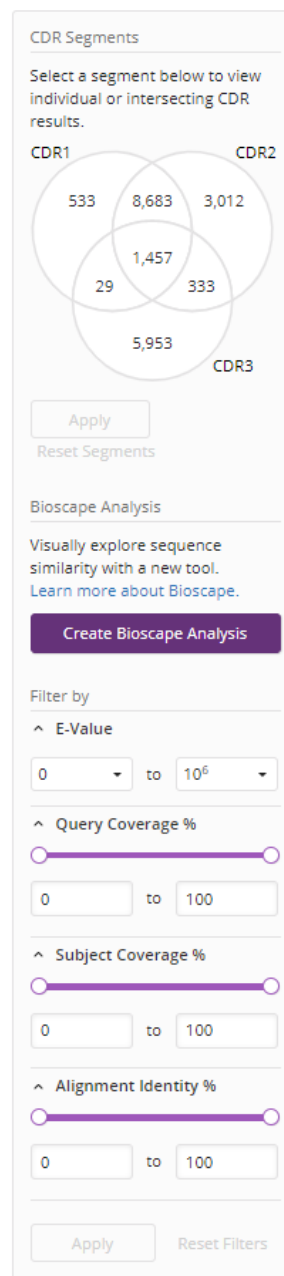
- Alignment visualization provides subject coverage, alignment position, and alignment identity at a glance

- Visually identify alignments to one, two, or all three of the CDR queries

# Motif Searching
## A sequence match search

**Biosequences**

Enter a protein or nucleotide string. Learn more about Biosequence Search.

| BLAST | CDR | Motif |

VGIGGGGGGGGGGXGGQGGXGX[SN][NG]XGGNGXGXGSHI

Clear Search

**Sequence Type:**

| Nucleotide | Protein |

Limit Total Sequence Results to:

20000 ▾

🔍 Start Biosequence Search

Advanced Biosequence Search ^          Reset All

Query Coverage % ❓     E-Value ❓

90          10 ▾

---

🧬 **Biosequences** (74)          View: Expanded ▾

📖 References                                      ⬇

Query Details    View More

> Seq 1: 1 VGIGGGGGGGGGGXGGQGGXGXSNXGGNGXGXGSHI 35 ▾

> **Seq 1: 1 VGIGGGGGGGGGGXGGQGGXGXSNXGGNGXGXGSHI 35**

> Seq 2: 1 VGIGGGGGGGGGGXGGQGGXGXSGXGGNGXGXGSHI 35

> Seq 3: 1 VGIGGGGGGGGGGXGGQGGXGXNNXGGNGXGXGSHI 35

> Seq 4: 1 VGIGGGGGGGGGGXGGQGGXGXNGXGGNGXGXGSHI 35

Matches: 35
Mismatches: 0

Subject ①————————————————————————————㉟

View Less ⌄

| Alignment | Subject | References |

📖 References

**Alignment Data**
BLAST Score: 266
E-Value: 1.26343e-35

Q        1  VGIGGGGGGG GGXGGQGGXG XSNXGGNGXG XGSHI  35
            |||||||||| |||||||||| |||||||||| |||||
S        1  VGIGGGGGGG GGXGGQGGXG XSNXGGNGXG XGSHI  35

---

[**SN**] in the above query finds sequences which contain either Serine (S) or Asparagine (N) at amino acid position 22.

X will find positive mismatches (+) for any hit amino acid or a match between query and hit X amino acid.

# Biosequence Motif Codes

| Degenerate Code | Meaning |
| --- | --- |
| X | Any amino acid |
| B | D or N |
| Z | E or Q |
| J | I or L |

| Degenerate Code | Meaning |
| --- | --- |
| N | A or C or G or T |
| R | A or G |
| Y | C or T |
| M | A or C |
| K | G or T |
| S | C or G |
| W | A or T |
| H | A or C or T |
| B | C or G or T |
| V | A or C or G |
| D | A or G or T |

CAS
A division of the
American Chemical Society

# Biosequence Motif Codes

| Degenerate Code | Meaning |
|---|---|
| [XYZ] | Any nucleotide or amino acid listed within the square brackets |
| {m,n} | At least m and maximum n residues length |
| {n} | Exactly n length |
| ^XYZ$ | Search for the exact sequence XYZ |

CAS

A division of the
American Chemical Society