# Biosequence Searching in

CAS STNext

**Access essential sequence databases for comprehensive searches**

- **CAS REGISTRY$^{SM}$**
- **DGENE**
- **PCTGEN**
- **USGENE$^{®}$**



**FIZ Karlsruhe**
Leibniz Institute for Information Infrastructure

**CAS**
A division of the
American Chemical Society

# Authority Sources in STNext

STNext allows you to combine and de-dup your results from the following sequence searchable databases:

CAS REGISTRY*

DGENE (Derwent Geneseq$^{TM}$)

PCTGEN

USGENE

By using a combination of these databases, you can locate even "hard-to-find" sequences. And, with STN you can gather them all into one, combined report complete with alignment and score matching details.



*Biosequence searching in the CAS REGISTRY database leverages the CAS Registry BLAST® client – which is a freely available utility.

# Before You Start Your Search

Save your sequence data in .TXT files. The sequence formats supported by STNext include: plain text, FASTA, GENBANK, and EMBL.

```
MSSPSLKWCF TLNYSSAAER ENFLSLLKEE DVHYAVVGDE VAPATGQKHL
QGYLSLKKRI RLGGLKKKYG SRAHWEIARG TDEENSKYCS KGTLILELGF
PVVNGSNKRK ISEMVARSPD RMKIEQPEIF HRYQSVNKLK KFKEEFVHPC
LDSPWQIQLT EAIDEEPDDR SIIWVYGPYG NEGKSTYAKS LIKKDWFYTR
```

```
>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas
maximus maximus]
LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATAFMGYVLPWGQMSFWGATV
ITNLFSAIPYIGTNLVEWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTF
```

```
  1 acaagatgcc attgtccccc ggcctcctgc tgctgctgct
 41 ctccggggcc acggccaccg ctgccctgcc cctggagggt
 81 ggccccaccg gccgagacag cgagcatatg caggaagcgg
121 caggaataag gaaaagcagc tcctgactt  tcctcgcttg
```

```
acaagatgcc attgtccccc ggcctcctgc tgctgctgct      40
ctccggggcc acggccaccg ctgccctgcc cctggagggt      80
ggccccaccg gccgagacag cgagcatatg caggaagcgg      120
caggaataag gaaaagcagc tcctgactt  tcctcgcttg      160
```

NOTE: Any spaces and numbers at the beginning or end of a line in the imported .TXT file will be stripped out when the sequence is uploaded into the session.

---

Prior to searching biosequences in CAS REGISTRY, download and install the CAS Registry BLAST client.

https://next.stn.org/stn/downloads/blast-download.html

NOTE: Your system administrator may need to assist you if you do not have installation privileges on your PC.

# Multi-database BLAST Search Strategy

Due to differences across databases, there are two different procedures for sequence searching in STNext. This example strategy walks through the steps of a comprehensive search that leverages all of the sequence searchable databases in STNext.

Example Search Scenario:
Find all **patents** disclosing the gene CBP1 from the soil bacterium *Serratia marcescens* with a minimum overall homology of 80%.

1. **Sequence Searching in DGENE, PCTGEN, or USGENE**

2. **Merge Answer Sets from DGENE, PCTGEN, USGENE**

3. **Search Using CAS Registry BLAST Client**

4. **Download Script and Alignment File**

5. **Run Script in STNext to Search CAplus$^{SM}$**

6. **Identify Duplicates**

7. **Create Combined Report**

# Sequence Searching in DGENE, PCTGEN, or USGENE

Login to STNext to get started. Use your STN login credentials at next.stn.org (Contact your STN helpdesk if you need assistance.)

1. File into DGENE.
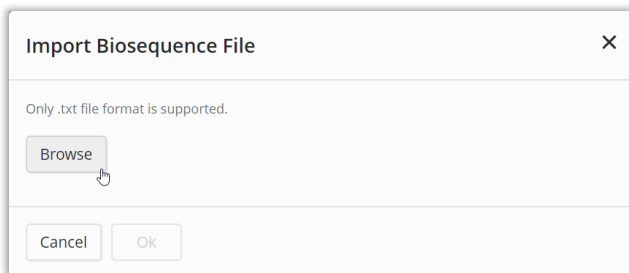
```
=> FIL DGENE
```

2. Select **Structures** from the My Files menu.



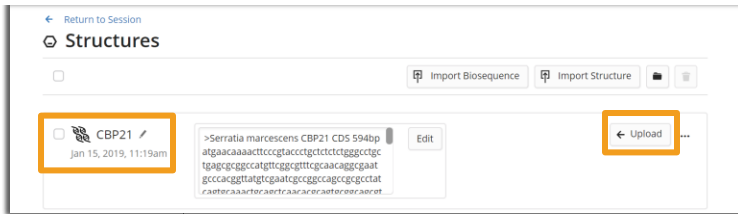3. Click the **Import Biosequence** button.



4. Browse to locate the .TXT file for the sequence



5. Click **OK** to import the file.

Sequences are indicated by the sequence icon and stored under My Files/Structures. Uploaded sequence queries may be up to 10,000 characters in length for BLAST search.



6. Click the **Upload** button.

   A sequence query L-number is automatically generated.

```
=> FIL DGENE
...

=>
Uploading sequence file: CBP21

UPLOAD SUCCESSFULLY COMPLETED
L1   GENERATED
```

(Option to verify your sequence using the D LQUE at this point.)

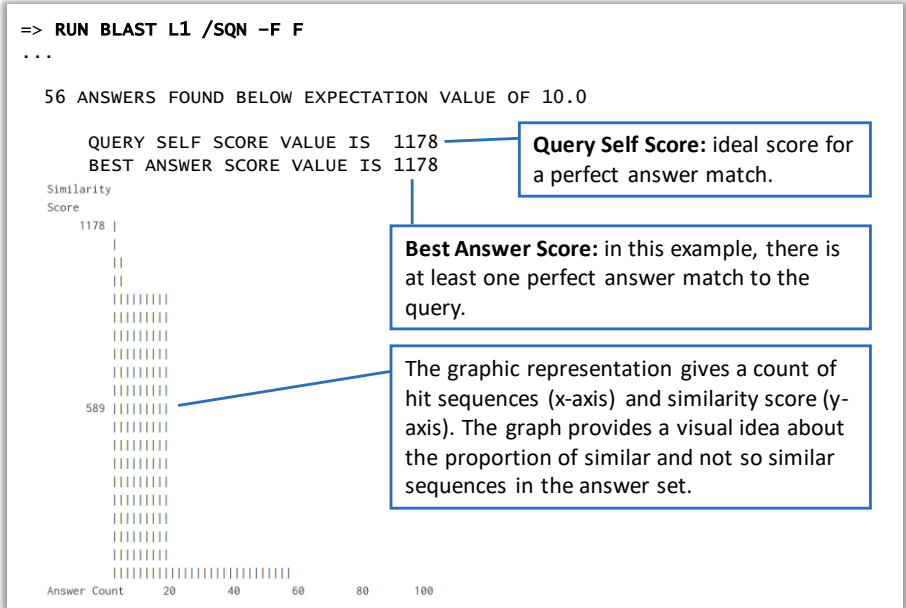7. Enter the BLAST command using the sequence query L-number.

```
=> RUN BLAST L1 /SQN –F F
```

The **low complexity filter** can eliminate biologically uninteresting segments that have low compositional complexity. The filter is set to F (False, off) by /SQN –F F (recommendation for patent sequence search).

| | |
|---|---|
| Protein search: | RUN BLAST L1 /**SQP** |
| Nucleotide search: | RUN BLAST L1 /**SQN** |
| Translated search: | RUN BLAST L1 /**TSQN** |

**BLAST:** NCBI BLAST for advanced similarity searching.

## 8. Evaluate the answer set.

```
=> RUN BLAST L1 /SQN -F F
...

   56 ANSWERS FOUND BELOW EXPECTATION VALUE OF 10.0

      QUERY SELF SCORE VALUE IS  1178
      BEST ANSWER SCORE VALUE IS 1178
```

**Query Self Score:** ideal score for a perfect answer match.

**Best Answer Score:** in this example, there is at least one perfect answer match to the query.

The graphic representation gives a count of hit sequences (x-axis) and similarity score (y-axis). The graph provides a visual idea about the proportion of similar and not so similar sequences in the answer set.

```
Similarity
Score
    1178 |
         |
         ||
         ||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
     589 |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||
         |||||||||||||||||||||||||||||||||||||
Answer Count    20    40    60    80   100
```

## 9. Decide how many answers to keep. You can choose a number of answers, all answers or a minimum percent self score. (This example uses 80% minimum score.)

```
ENTER EITHER THE NUMBER OF ANSWERS YOU WISH TO KEEP
 OR ENTER MINIMUM PERCENT OF SELF SCORE FOLLOWED BY %
 (BEST ANSWER PERCENTAGE OF SELF SCORE IS 100%)
ENTER (ALL) OR ? :  80%

L2    RUN STATEMENT CREATED
L2         56 ATGAACAAAACTTCCCGTACCCTGCTCTCTCTGGGCCTGCTGAGCGCGGC
              CATGTTCGGCGTTTCGCAACAGGCGAATGCCCACGGTTATGTCGAATCGC
              CGGCCAGCCGCGCCTATCAGTGCAAACTGCAGCTCAACACGCAGTGCGGC
              AGCGTGCAGTACGAACCGCAGAGCGTCGAGGGCCTGAAAGGCTTCCCGCA
              GGCCGGCCCGGCTGACGGCCATATCGCCAGCGCCGACAAGTCCACCTTCT
              TCGAACTGGATCAGCAAACGCCGACGCGCTGGAA/SQN -F F
```

## 10. Use a Display command to review sample records.

```
=> D TRIAL SCORE ALIGN 1-20
```

## 11. Repeat this sequence search in PCTGEN and USGENE.

# Merge Answer Sets from DGENE, PCTGEN, USGENE

When you have L-numbers for the DGENE, PCTGEN and USGENE sequence search results, merge those results into a single answer set.

1. Enter the SET DUPORDER FILE setting command to specify the answer retrieval follows the order the files were searched.

```
=> SET DUPORDER FILE

SET COMMAND COMPLETED
```

2. Enter the DUPLICATE IDENTIFY (DUP IDE) command to create a new L-number from the multiple answer sets.

```
=> DUP IDE L2 L3 L4

FILE 'DGENE' ENTERED AT 13:06:37 ON 17 JAN 2019
COPYRIGHT (C) 2019 CLARIVATE ANALYTICS
FILE 'USGENE' ENTERED AT 13:06:37 ON 17 JAN 2019
COPYRIGHT (C) 2019 SEQUENCEBASE CORP
FILE 'PCTGEN' ENTERED AT 13:06:37 ON 17 JAN 2019
COPYRIGHT (C) 2019 WIPO

PROCESSING COMPLETED FOR L2
PROCESSING COMPLETED FOR L3
PROCESSING COMPLETED FOR L4
L5          132 DUP IDE L2 L3 L4 (INCLUDES 0 SETS OF DUPLICATES)
                ANSWERS '1-56' FROM FILE DGENE
                ANSWERS '57-121' FROM FILE USGENE
                ANSWERS '122-132' FROM FILE PCTGEN
```

Leveraging the DUP IDE command in this way is simply getting all answers into a single answer set.

3. Sort the results by descending similarity score and identity using the SOR SCORE D IDENT D command. Sorting is useful when there are many answers to review.

```
=> SOR L5 SCORE D IDENT D

PROCESSING COMPLETED FOR L6
L6          132 SOR L5 SCORE D IDENT D
```

4. Click the ellipsis (. . .) button and select the Patent Family Manager to display the results. (This example specifically targets patents.)



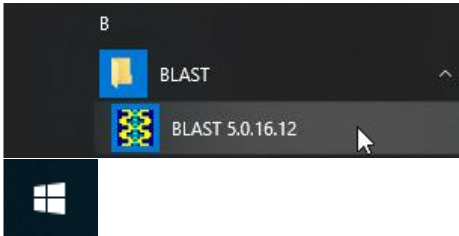5. Choose a patent family display option.



This example uses BIB SQL SCORE IDENT ALIGN for the First Member of Each Family.
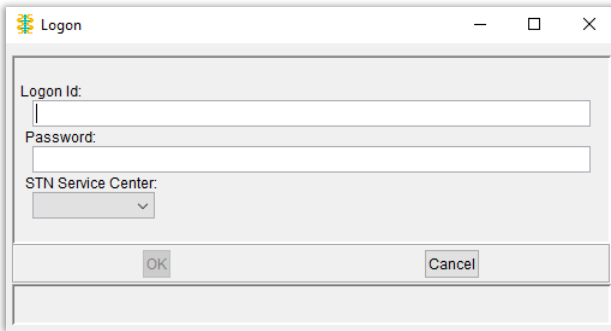
6. Click the **Submit** button.

STNext's Patent Family Manager automatically starts with an FSORT to place records in extended patent families – generating a new L-number and then proceeding with the record display.

# Search Using CAS Registry BLAST Client

After installing the CAS Registry BLAST client on your PC, locate the link/ icon under the Start menu.



1. Click the **BLAST** icon.

2. Use your STN login credentials and choose your local STN Service Center from the drop-down. Click **OK** to login.



3. Click the **New Search** button.

4. Click the **Similar Sequences** button.



5. Name the result and enter the sequence by copy/paste or reading data from a file.



6. Click the **OK** button.

7. Click the appropriate BLAST query type. (This example uses BLASTn.)



8. Choose the nucleotide subsets desired. (This example uses All.)

9. Click the **OK** button to save the subset selection.

10. Adjust the BLAST settings. (This example is patent specific, so uncheck the Low Complexity Filtering option.)



11. Click the **OK** button to start the search.

The Result Set Manager will display the sequence search and status. You can submit additional sequences while the searches are running. Up to 100 results sets can be kept in the Results Set Manager.

When a report is complete, you can view results.

# Download Script and Alignment File from CAS Registry

When your BLAST query is complete, you can evaluate the results and select which answers to download for use in STNext.



1. Select the results set from the **Reports** tab and click the **View Results** button.

2. Review key sequence statistics such as total unique and redundant sequences. There is also a grouping of results based on alignment scores.


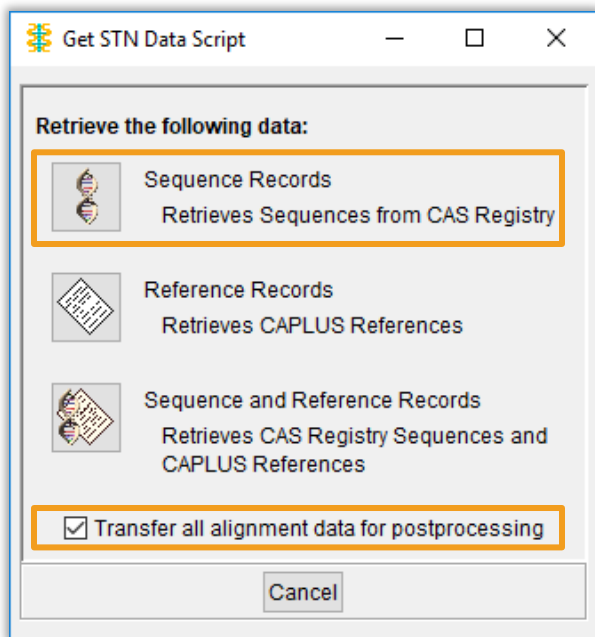
Higher Alignment Scores will have better alignment and match over the length of the query.

3. Click the **Alignment Score** buttons to select groups of sequences to include in the STN Data Script.

4. Preview the alignment detail for individual results to verify the letter by letter alignment details of the selected results, if desired.

5. Click the **Get STN Data Script** button.

6. Check the **Transfer all alignment data for postprocessing** option.



7. Click the **Sequence Records** button.

The system will prompt you to name and save a .SCB script file that will be used in STNext.

The system will then prompt you a second time to save a .XSS file that contains the alignment data that will be available for use in your combined report.

Note the file directory so it is easy to locate and import the .SCB and .XSS files into STNext.

## Run Script in STNext to Search CAplus

1. Select **Scripts** from the My Files menu.



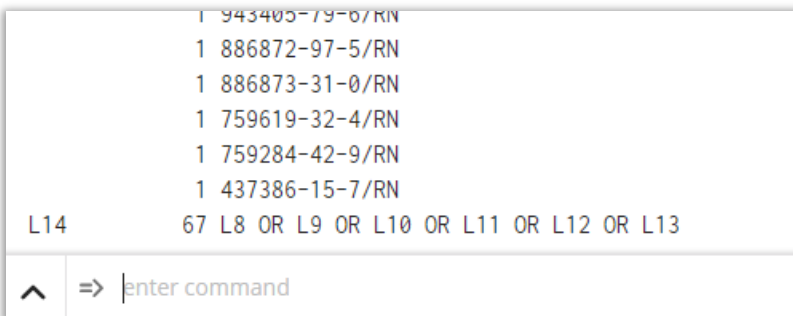2. Click the **Import Script** button.



3. Browse to locate the .SCB script file that was created in the CAS Registry BLAST Client. Click the **OK** button.



4. The .SCB file is saved to the Scripts page. Click the **Run** button.



The script will automatically file into CAS REGISTRY and search for the CAS RNs selected in the BLAST client session and combine the RNs in a single L-number.

## Identify Duplicates

When the .SCB script runs, it finds *SUBSTANCE* records in CAS REGISTRY related to the sequences. To combine these answers with those found in DGENE, PCTGEN or USGENE, we have to retrieve the corresponding *DOCUMENT* records. Then we can compare all the answers and avoid duplicates.

1. File into CAplus or HCAplus.

```
...
L14          67 L8 OR L9 OR L10 OR L11 OR L12 OR L13

=> FIL CAPLUS
```

2. Search L14 for corresponding patent records by adding the P/DT command. (Patent/Document Type)

```
=> S L14 AND P/DT

L15          19 L14 AND P/DT
```

3. Transfer patent numbers from the DGENE, PCTGEN and USGENE answers to CAplus.

```
=> TRA L6 1- PN

L16        TRANSFER L6 1- PN :     39 TERMS
L17        38 L16
L18          QUE  TERMS FROM L16 WITH NO HITS:    1 TERM
```

18

4. Search both the CAplus and DGENE/PCTGEN/USGENE answer sets using the NOT operator to identify the unique records from the CAplus search.

```
=> S L15 NOT L17

L19              9 L15 NOT L17
```

In this example, there were 9 additional records found.


5. Display those unique records using the HITRN command so the CAS Registry BLAST alignments can be included in the combined report.
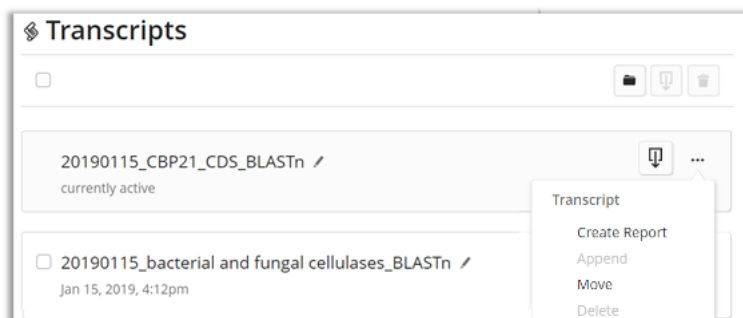
```
=> D L19 1- BIB HITRN

L19  ANSWER 1 OF 9  HCAPLUS  COPYRIGHT 2019 ACS on STN
PatentPak PDF | PatentPak PDF+ | PatentPak Interactive
AN   2018:864066  HCAPLUS Full-text
DN   168:509520
TI   Highly active self-sufficient nitration biocatalysts based on chimeric
     cytochrome P 450 enzymes fusion proteins
IN   Ding, Yousong; Zou, Ran
PA   University of Florida Research Foundation, Incorporated, USA
SO   PCT Int. Appl., 221pp.
     CODEN: PIXXD2
DT   Patent
LA   English
FAN.CNT 1
PPPI
     PATENT NO.         KIND  DATE      LANGUAGE   PatentPak
     WO 2018081456      A1    20180503  English    PDF | PDF+ | Interactive
...
IT   2222417-62-9
     RL: BSU (Biological study, unclassified); PRP (Properties); BIOL
     (Biological study)
        (nucleotide sequence; highly active self-sufficient nitration
        biocatalysts based on chimeric cytochrome P 450 enzymes fusion
        proteins)
```
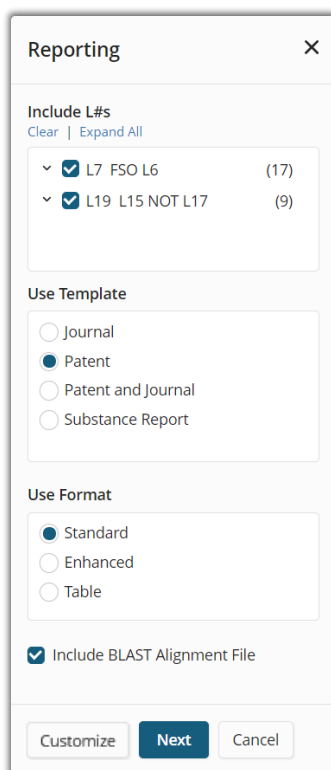
# Create Combined Report

After collecting answers from all the databases, we can create a single report that also includes the alignment data.
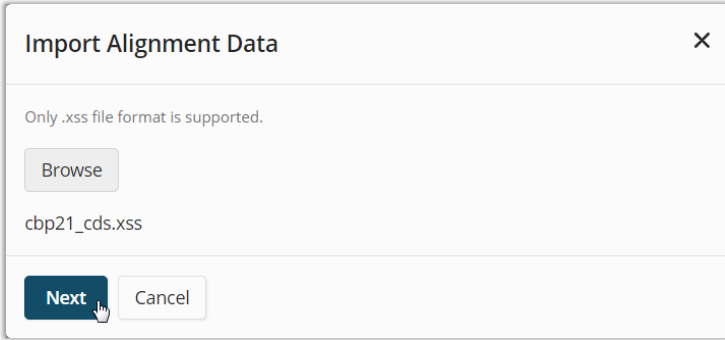
1. Select **Transcripts** from the My Files menu.

2. Click the ellipsis (. . .) button and select **Create Report**.



3. Select the L7 records sorted by patent family, and L19 records that were unique to the CAplus search.

4. This example scenario was patent-specific, so use the **Patent** template and **Standard** format.

5. Check the **Include BLAST Alignment File** option.

6. Click the **Next** button.

7. Browse to locate the .XSS alignment file. Click the **Next** button.

**Import Alignment Data**   ✕

Only .xss file format is supported.

Browse

cbp21_cds.xss

[ Next ]  [ Cancel ]

8. Enter Report Header information as desired.

**Document Header**   ✕

Document Title

CBP21_BLASTn

Creator

E Aichinger

Comments

S marcescens CBP21 cds
DGENE PCTGEN USGENE CAS Reg BLAST
all wo cut off

Include:
- ☐ Page Numbers
- ☐ Date and Time
- ☐ Cover Page

[ Download ]  [ Cancel ]

9. Click the **Download** button.

# Summary

STNext offers four essential patent sequence databases with unrivaled coverage for comprehensive searching.

The cross-over capability allows you to combine data from several sources into one report that includes query information, record details and BLAST alignment data.

Use your STN login credentials to access STNext at
next.stn.org

**Record Details**

**BLAST Alignment Data**

# Download the CAS Registry BLAST® client

Prior to searching biosequences in CAS REGISTRY, download and install the CAS Registry BLAST client.

With the CAS REGISTRY$^{SM}$ database, the CAS Registry BLAST® client must be used to locate the sequence data as CAS Registry Numbers® and export a script that will help load the Registry Numbers into your session on STNext.

Download the CAS Registry BLAST client at:
https://next.stn.org/stn/downloads/blast-download.html

Please note, your system administrator may need to assist you if you do not have installation privileges on your PC.

## For more information...

CAS
help@cas.org

Support & Training:
www.cas.org/support

FIZ Karlsruhe
helpdesk@fiz-karlsruhe.de

Support & Training:
www.stn-international.de