

# STN<sup>®</sup>

Resources for Intellectual Property  
Sequence searching on STN<sup>®</sup>

Robert Austin – FIZ Karlsruhe

## Agenda

2

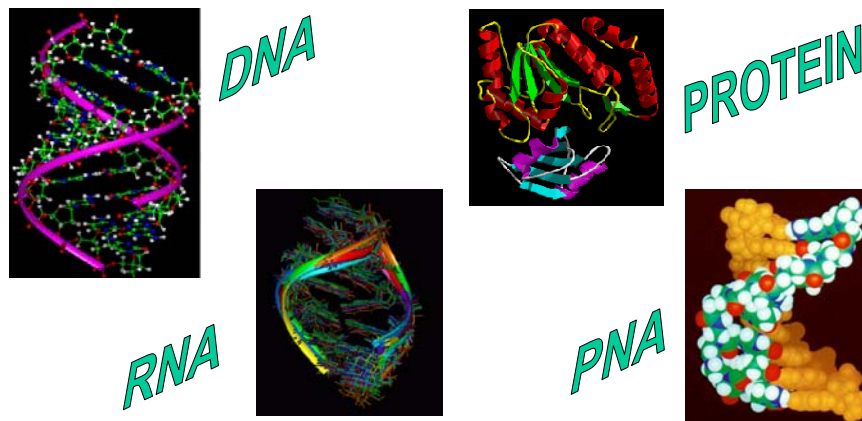
- Sequences in patent publications
- STN patent sequence databases
- Comparison to web based resources
- Results of a comparative search example
- Summary and resources

**STN**

 FIZ Karlsruhe

## A completely new subject for many....

3



Bluff Your Way in Genetics!!

[http://www.stn-international.com/training\\_center/bioseq/bluff.pdf](http://www.stn-international.com/training_center/bioseq/bluff.pdf)

**STN**

FIZ Karlsruhe

## Sequences are often described in patents, and usually via a sequence identity number

4

```
L1 ANSWER 1 OF 1 USGENE COPYRIGHT 2008 SEQUENCEBASE CORP on STN
CLM US6831165 B1: What is claimed is:
1. An isolated nucleic acid molecule comprising the nucleotide
sequence of SEQ ID NO:1, or the full complement thereof.

2. An isolated nucleic acid molecule which encodes a polypeptide
comprising the amino acid sequence of SEQ ID NO:2.

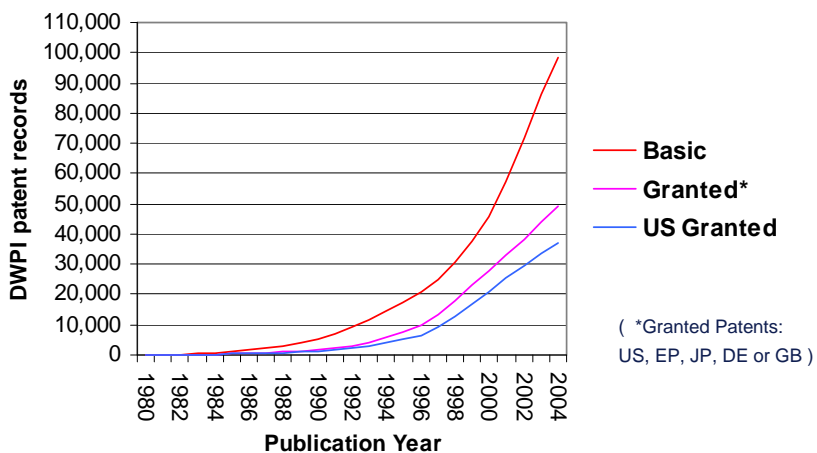
3. An isolated nucleic acid molecule which encodes a polypeptide
comprising the amino acid sequence of SEQ ID NO:2, wherein the
nucleic acid molecule hybridizes to the complement of a nucleic
acid molecule consisting of SEQ ID NO:1 in 6.times.SSC at 45.degree.
C., followed by one or more washes in 0.2.times.SSC, 0.1% SDS at 50-
65.degree. C., and wherein said nucleic acid molecule encodes a
sulfate adenylate transferase subunit 2.
. . . .
```

**STN**

FIZ Karlsruhe

## Cumulative growth of patent publications with sequences

5



STN

FIZ Karlsruhe

## Sequences referred to by SEQ ID NO typically appear in a Sequence Listing

6

```

<210> SEQ ID NO 1
<211> LENGTH: 314
<212> TYPE: DNA
<213> ORGANISM: Corynebacterium glutamicum
<220> FEATURE:
<221> NAME/KEY: CDS
<222> LOCATION: (1)..(291)
<223> OTHER INFORMATION: RXA02548
<400> SEQUENCE: 1
  cca ccg atc tac ttc tcc cac gac cgc gaa gtt ttc gag cgc gac ggc   48
  Pro Pro Ile Tyr Phe Ser His Asp Arg Glu Val Phe Glu Arg Asp Gly
   1           5           10          15
  atg tgg ctg acc gca ggc gag tgg ggt gga cca aag aag ggc gag gag   96
  Met Trp Leu Thr Ala Gly Glu Trp Gly Gly Pro Lys Lys Gly Glu Glu
           20           25           30
  atc gtc acc aag act gtc cgc tac cgc acc gtc ggc gat atg tcc tgc  144
  Ile Val Thr
           35
  acc ggt gct
  . . . .
    
```

STN databases, DGENE, USGENE, PCTGEN and REGISTRY take patent sequence data like this and provide it in a suitably searchable format.

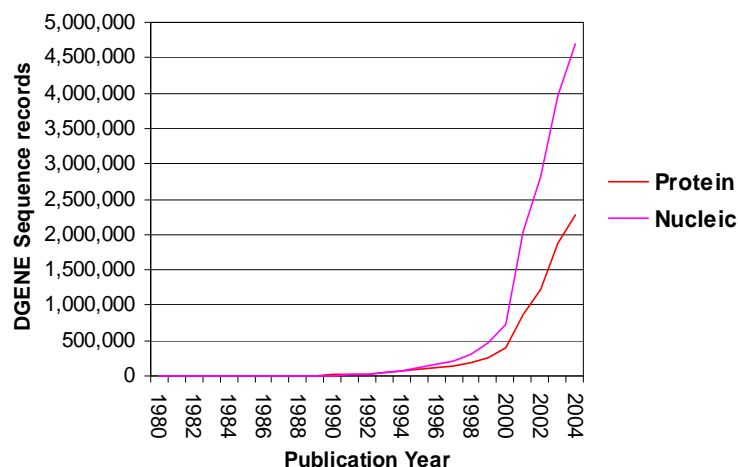
For further details about WIPO ST.25 format, visit:  
<http://www.wipo.int/scit/en/standards/pdf/03-25-01.pdf>

STN

FIZ Karlsruhe

## Cumulative growth of sequences published in patent literature

7



**STN**

**FIZ Karlsruhe**

## Agenda

8

- Sequences in patent publications
- **STN patent sequence databases**
- Comparison to web based resources
- Results of a comparative search example
- Summary and resources

**STN**

**FIZ Karlsruhe**

## STN sequence searchable databases

9

- DGENE
  - Thomson Reuters GENESEQ™
  - Value-added patent sequence data from around the globe
- USGENE®
  - The USPTO Genetic Sequence Database
  - A new and unique access point to USPTO sequence data
- PCTGEN
  - WIPO/PCT Patent Application Biosequences
  - The complete collection of e-published sequences from WIPO
- CAS REGISTRY<sup>SM</sup>
  - Chemical Abstracts Service (CAS) REGISTRY File
  - Worldwide value-added patent and non-patent sequence data

**STN**

 FIZ Karlsruhe

## Thomson Reuters GENESEQ (DGENE)

10

- Largest value-added patent sequence database
- Used routinely by all major patent offices\*
- Sequences from the basic patents of the 40 authorities of the *Derwent World Patents Index*®
- Bibliography, enhanced title, abstract, indexing and patent location provided for each sequence
- Patent Family and Legal Status display
- Updated every two weeks
- 1981 - present

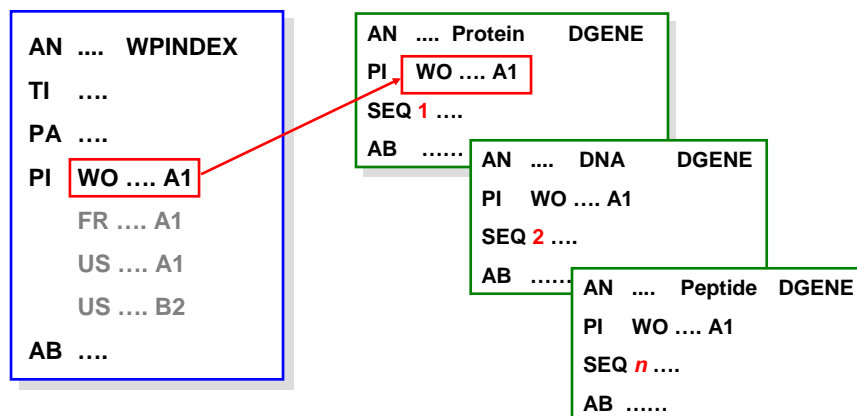
\* See page 11: [http://www.trilateral.net/projects/biotechnology/search\\_guidebook\\_vers\\_1.pdf](http://www.trilateral.net/projects/biotechnology/search_guidebook_vers_1.pdf)

**STN**

 FIZ Karlsruhe

## Relationship between DWPI patent family and DGENE sequence database

11



WPINDEX = Derwent World Patents Index® (DWPI<sup>SM</sup>) on STN®

DGENE = GENESEQ<sup>TM</sup> on STN®

**STN**

FIZ Karlsruhe

## What exactly is the “value-add” in DGENE?

12

- DWPI patent title, concise sequence description, abstract and keyword indexing *per sequence*
  - Illuminate the context of *each sequence* within the invention
  - Superior text based refinement of sequence searches
  - Efficient scanning and review of search results for relevance
- Feature tables for sequence modifications/annotations
  - Extensive detailed annotations are provided by Indexers
- Patent sequence location (claim, example, etc)
  - Assigned manually by Thomson Indexers
  - Ability to filter searches to those described in the claims
- Sequences intellectually derived by Indexers
  - Provides unique sequence hits not disclosed in formal listings

**STN**

FIZ Karlsruhe

## Some editorial insights regarding WIPO/PCT sequences indexed in DGENE 13

- On average 120 WIPO/PCT basic patents have sequences indexed into DGENE each week
- Of these, about 15-20 may have electronic listings available – the rest are keyed manually
  - Sequences are independently double-keyed with a guaranteed accuracy of 99.995% (1 in 20,000)
- About 15% of PCTs with electronic listings have extra sequences indexed from the specification
- Typically 1 or 2 documents per week will also have intellectually derived sequences indexed, based upon the wording of the patent claims

Source: Colin Williams, GENESEQ Editorial & Content Manager, Thomson Reuters (12/2006)



## Derived sequences are intellectually created by Indexers from wording in the patent text 14

AN	AEJ92622	protein	DGENE
TI	Hydrolyzing/synthesizing carboxylic acid ester/amide from chiral/prochiral reactants for preparing e.g. pharmaceuticals, comprises contacting reactants with a polypeptide having hydrolytic activity.		
IN	Svendsen A; Vind J; Brask J; De Maria L		
PA	(NOVO)	NOVOZYMES AS.	
PI	WO 2006084470	A2	20060817
AI	WO 2006-DK76		20060210
PRAI	EP 2005-388012		20050210
PSL	Claim 16		
DED	19 OCT 2006	(first entry)	
LA	English		
OS	2006-560037 [57]		
DESC	Variant fungal lipolytic hydrolase #2.		
KW	hydrolysis; lipase; pharmaceutical; pesticide; enzyme; mutein.		
ORGN	Thermomyces lanuginosus. Synthetic.		
AB	The new invention relates to a enzymatic method of hydrolyzing or synthesizing carboxylic acid ester or amide from chiral or prochiral reactants, by providing reactants for hydrolysis or synthesis, and contacting the reactants with a polypeptide which has hydrolytic activity on ester or amide, and a sequence 50% homologous to Thermomyces lanuginosus lipase. Also described is a polypeptide, which has hydrolase activity on an ester or amide substrate, and has an amino acid sequence that has at least 80% identity to SEQ ID No: 5 and compared to SEQ ID No: 5 comprises a substitution corresponding to I90Q, N92TD, F95Y, F113Y, I202M, V203GM, L269T and 270F. . . . .		

In this example, the indexer has intellectually derived this sequence from the wild-type lipolytic hydrolase.



## Indexers explain exactly how they derived the sequence at the end of the abstract

15

The polypeptide is at least 80% homologous to any of SEQ ID No: 1-6 being amino acid sequences of lipolytic enzymes of fungus such as *Rhizomucor miehei* (SWISSPROT P19515), *Rhizopus delemar*, *Fusarium oxysporum*, *Penicillium camemberti* (SWISSPROT P25234), *Thermomyces lanuginosus* (SWISSPROT 059952) and *Thermomyces ibadanensis* . . . . . The method is useful in the preparation of pharmaceuticals or pesticides, where the synthesis includes synthesis of 2-butyl propionate. This sequence is a variant fungal lipolytic hydrolase (lipase), V203M T231R N233R. This sequence is not shown in the specification, but was created by the indexer using the information given in claim 16.

```
SQL 269
SEQ 1 evsqdlfnqf nlfaqysaaa ycgknndapa gtnitctgna cpevekadat
51 flysfedsgv gdtvgflald ntnklivlsf rgsrsienwi glnlnfdlkei
101 ndicsgcrgh dgftsswrsv adtlrpkved avrehpdyrv vftghslgga
151 latvagadlr gng
201 dimprlppre fgy
251 nipdipahlw yfg
```

The indexer has added explanatory sentences to the abstract, and annotations to the feature table.

### FEATURE TABLE:

Key	Location	Qualifier	
Modified-site	203	note	"Wild type Val replaced by Met"
Modified-site	231	note	"Wild type Thr replaced by Arg"
Modified-site	233	note	"Wild type Asn replaced by Arg"

STN

FIZ Karlsruhe

## USGENE is the USPTO Genetic Sequence Database

16

- Sequences from all relevant USPTO published patent applications and granted (issued) patents
- Assignee and full inventor names; publication, application and parent case PCT numbers and dates; original publication **title**, **abstract**, and **claims**
- Organism name, sequence length, Molecule Type, SEQ ID, and feature tables for features/annotations
- Produced by the SequenceBase Corporation
- Updated weekly – within **3 days** of publication
- 1982 – present

STN

FIZ Karlsruhe



## An individual publication is represented by one or more USGENE sequence records

17

(12) <b>United States Patent</b> Cutitta et al.	(10) Patent No.: <b>US 7,364,719 B2</b> (45) Date of Patent: <b>Apr. 29, 2008</b>	<b>AN ... Protein USGENE</b> <b>PI US ... B2</b> <b>SEQ 1 ...</b>
(54) <b>VASOREGULATING COMPOUNDS AND METHODS OF THEIR USE.</b>	6,440,421 B1 8/2002 Cornish et al. 2002/0055615 A1 5/2002 Cutitta et al.	<b>AN ... DNA USGENE</b> <b>PI US ... B2</b> <b>SEQ 2 ...</b>
(75) Inventors: <b>Frank Cutitta</b> , Adamstown, MD (US); <b>Alfredo Martinez</b> , Bethesda, MD (US); <b>William G. Stetler-Stevenson</b> , Kensington, MD (US); <b>Edward J.</b> <b>Uusivirta</b> , Kensington, MD (US); <b>Juan</b> <b>M. Saavedra</b> , Bethesda, MD (US)	FOREIGN PATENT DOCUMENTS EP 0 845 036 6/1999 EP 0 926 238 A2 11/2000 EP 0 926 238 A3 11/2000 WO 97/07214 2/1997 WO 01/18550 3/2001 WO 2004/043383 5/2004	<b>AN ... cDNA USGENE</b> <b>PI US ... B2</b> <b>SEQ n ...</b>
(73) Assignee: <b>The United States of America as represented by the Department of Health and Human Services.</b> Washington, DC (US)	OTHER PUBLICATIONS Corti et al., "Vasopressin Inhibitors: A New Therapeutic Concept in Cardiovascular Disease," <i>Cardiovascular Drugs</i> 104:1856-1862 (Oct. 9, 2001). Fernandez-Paton, "Vascular Matrix Metalloproteinase-2-Dependent Clearance of Calcitonin Gene-Related Peptide Promotes Vasodilation," <i>Circ. Res.</i> 87:676-679 (2000). Kizumura et al., "Cloning and characterization of cDNA encoding a precursor for human adrenomedullin," <i>Biochem. Biophys. Res. Comm.</i> 194:720-725 (1993). Kizumura et al., "Adrenomedullin (11-26): a novel endogenous hypotensive peptide isolated from bovine adrenal medulla," <i>Peptides</i> 22:1713-1718 (2001). Lewis et al., "Degradation of human adrenomedullin (1-52) by plasma membrane enzymes and identification of metabolites," <i>Peptides</i> 18(5):733-739 (1997). Watanabe et al., "Vasopressin activity of N-terminal fragments of adrenomedullin in anesthetized rat," <i>Biochem. Biophys. Res. Comm.</i> 219:643 (1996). Belloni et al., "Proadrenomedullin N-Terminal 20 Peptide (PAMP), Acting Through PAMH2(20)-Sensitive Receptors, Initiates Ca <sup>2+</sup> -Dependent, Agonist-Stimulated Secretion of Human Adrenal Glands," <i>Hypertension</i> 33:1185-1189 (1999). Calvo et al., "Adrenomedullin and proadrenomedullin N-terminal 20 peptide in the normal prostate and in prostate carcinoma," <i>Hormone. Res.</i> 57(2):98-104 (Apr. 2002) <i>Abstract Only</i> . Champion et al., "Proadrenomedullin N12-terminal 20 peptide has direct vasodilator activity in the rat," <i>Am. J. Physiol.</i> 272(4 Pt 2):R1047-54 (Apr. 1997) <i>Abstract Only</i> . Champion et al., "Ions-dependent vasodilator responses to proadrenomedullin N12-terminal 20 peptide in the hindquarters vascular bed of the rat," <i>Peptides</i> 18(4):513-519 (1997) <i>Abstract Only</i> .	
(*) Notice: Subject to any disclaimer, the term of this patent is extended or adjusted under 35 U.S.C. 154(b) by 179 days.		
(21) Appl. No.: <b>10/529,118</b>		
(22) PCT Filed: <b>Oct. 3, 2003</b>		
(86) PCT No.: <b>PCT/US03/31400</b> § 371 (c)(1), (2), (4) Date: <b>Mar. 24, 2005</b>		
(87) PCT Pub. No.: <b>WO/2004/032708</b> PCT Pub. Date: <b>Apr. 22, 2004</b>		
(65) <b>Prior Publication Data</b> US 2005/0261179 A1 Nov. 24, 2005		
<b>Related U.S. Application Data</b> (60) Provisional application No. 60/416,291, filed on Oct. 4, 2002.		
(51) Int. Cl. <b>A61K 49/00</b> (2006.01)		

STN

FIZ Karlsruhe

## USGENE consolidates unique USPTO sequence data from different sources

18

- USPTO Publication Site for Issued and Published Sequences (PSIPS)
  - The official mega-publication download site, 2001-date
- International Nucleotide Sequence Database Collaboration (INSDC) (NCBI/EMBL/DDBJ, Genbank)
  - U.S. granted patent nucleotide sequences, 1982-date
- USPTO Protein Database (NCBI/EMBL)
  - U.S. granted patent protein/peptide sequences, 1982-date
- USPTO Published Applications and Patents Full-Text
  - Filling in omissions, coverage gaps and to enhance timeliness

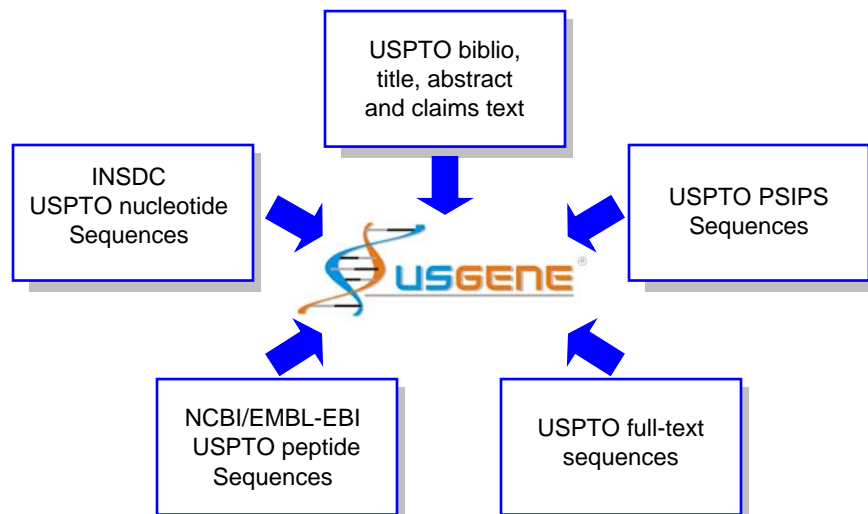
The USGENE Sequence Source (/SSO) field indicates which source any given USGENE sequence record was derived from.

STN

FIZ Karlsruhe

## USGENE combines these sequences with bibliographic data and claims text

19



STN

FIZ Karlsruhe

## USGENE is an essential additional tool for tackling business critical searches

20

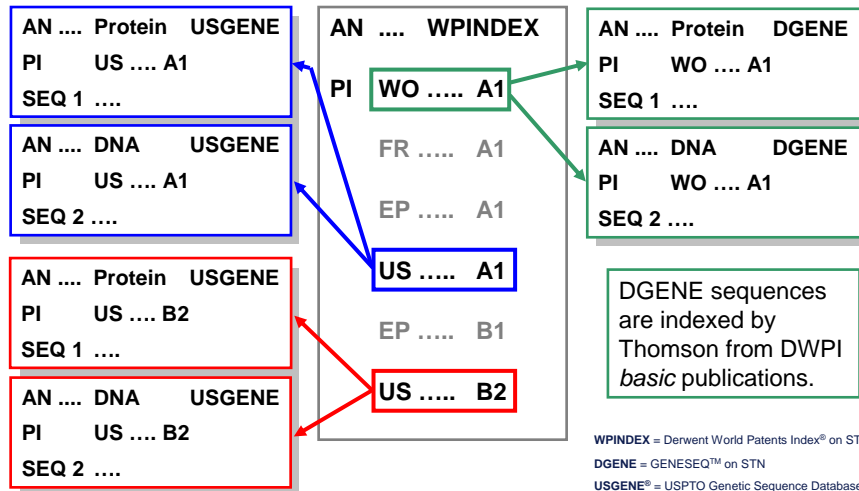
- DGENE provides curated and indexed patent sequence data from the DWPI *basic* publication
  - 61% of *basics* are WIPO/PCT published applications
  - Updated biweekly, typically 65 days from publication
- USGENE provides all available sequence data from the USPTO as a single merged resource
  - Both **U.S. patents** and **U.S. published applications**
  - Updated weekly, within **3 days** of USPTO publication
- Sequence listing variation often occurs between PCT and U.S. granted patent publication stages
  - Especially important, e.g. for freedom-to-operate

STN

FIZ Karlsruhe

## USGENE and DGENE capture sequence data from different patent family members

21



STN

FIZ Karlsruhe

## Sequence listing variation often occurs between PCT and U.S. granted patent stage

22

```

L1 ANSWER 1 OF 1 WPINDEX COPYRIGHT 2008 THOMSON REUTERS on STN
AN 1994-358278 [44] WPINDEX
TI New polynucleotide(s) specific for hepatitis C virus types 4, 5 and 6 -
and related antigenic peptide(s) and antibodies, useful in vaccines,
diagnosis, HCV typing and treatment
DC B04; D16; S03
IN PIKE I H; SIMMONDS P; YAP P L
PA (COMM-N) COMMON SERVICES AGENCY; (MURE-N) MUREX DIAGNOSTICS INT INC; . . .
PI WO 9425602 A1 19941110 (199444)* EN 70[5]
AU 9465797 A 19941110 (199444)* EN 70[5]
FI 9505224 A 19941110 (199444)* EN 70[5]
EP 698101 A1 19941110 (199444)* EN 70[5]
JP 09500009 W 19941110 (199444)* EN 70[5]
AU 695259 B 19941110 (199444)* EN 70[5]
EP 698101 B1 20050210 (200512) EN
DE 69434116 E 20050419 (200527) EN
US 20050032047 A1 20050210 (200512) EN
US 6881821 B2 20050419 (200527) EN
. . . . .
ADT WO 9425602 A1 WO 1994-GB957 19940505 . . . .
PRAT GB 1994-263 19940107
GB 1993-9237 19930505
    
```

In this example the patent family has:

- 9 sequences from WO9425602 in DGENE
- 50 sequences from US20050032047 in USGENE
- 58 sequences from US6881821 in USGENE

STN

FIZ Karlsruhe

## PCTGEN is the World Patent Application Biosequences database on STN

23

- Produced by FIZ Karlsruhe and WIPO
- Sequences submitted & published electronically as a formal part of PCT patent applications
- Publication number and date, patent applicant name(s) and the original publication title are provided for each sequence
- Sequence length, SEQ ID, organism name and molecule type are included for each sequence
- Updated weekly – within **24 hours** of publication
- August 2001 – present

**STN**

FIZ Karlsruhe

## Relationship between PCTFULL and PCTGEN databases

24

AN ... PCTFULL

TI ....

PA ....

PI WO .... A1

AB ....

DETD ....

CLM ....

AN ... Protein PCTGEN

PI WO .... A1

SEQ 1 ....

AN ... DNA PCTGEN

PI WO .... A1

SEQ 2 ....

AN ... Peptide PCTGEN

PI WO .... A1

SEQ *n* ....

PCTFULL = WIPO/PCT patent applications full-text

PCTGEN = WIPO/PCT patent application biosequences

**STN**

FIZ Karlsruhe

## Chemical Abstracts Service (CAS) REGISTRY database

25

- Sequences from >3000 life science journals
- Sequences from the basic patents of the 50 patent authorities of the CAPLUS file on STN
- Patent number, location and standardized nomenclature are provided for each sequence
- Updated daily
- 1907 - present

**STN**

FIZ Karlsruhe

## Relationship between CAPLUS patent family and CAS Registry databases

26

AN .... CAPLUS

TI ....

PA ....

PI WO .... A1

FR .... A1

US .... A1

US .... B2

AB ....

IT RN ....

RN .... Protein REGISTRY

PI WO .... A1

SEQ 1 ....

RN .... DNA REGISTRY

PI WO .... A1

SEQ 2 ....

RN .... Peptide REGISTRY

PI WO .... A1

SEQ *n* ....

**STN**

FIZ Karlsruhe

## Agenda

27

- Sequences in patent publications
- STN patent sequence databases
- **Comparison to web based resources**
- Results of a comparative search example
- Summary and resources

**STN**

 FIZ Karlsruhe

## There are three main web resources which provide searchable patent sequence data

28

- National Center for Biotechnology Information (NCBI) of the U.S. National Library of Medicine
  - [www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)
- European Bioinformatics Institute\* (EBI)
  - [www.ebi.ac.uk](http://www.ebi.ac.uk)
- DNA DataBank of Japan (DDBJ)
  - [www.ddbj.nig.ac.jp](http://www.ddbj.nig.ac.jp)
- The USPTO, EPO and JPO rely on the NCBI, EBI and DDBJ, respectively, to provide an interface for searching patent sequence data

(\* The EBI is the U.K. based outstation of the European Molecular Biology Laboratory – EMBL)

**STN**

 FIZ Karlsruhe

## STN databases have more patent sequences than web resources

29

<u>Database</u>	<u>Peptide</u>	<u>Nucleotide</u>
REGISTRY	2.50	9.39
DGENE	2.71	5.47
PCTGEN	0.54	4.01
NCBI	0.57	3.43

(Statistics in millions: October 8<sup>th</sup>, 2006)



## Why does STN have more patent sequences than the NCBI?

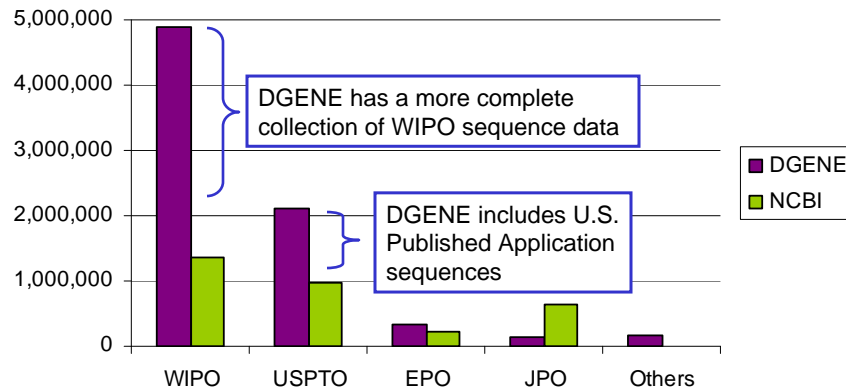
30

1. NCBI has a substantially incomplete collection of WIPO/PCT sequence data
2. NCBI does not cover any USPTO Published Application sequence data
3. DGENE and REGISTRY cover many more patent authorities than NCBI
4. DGENE, USGENE, PCTGEN and REGISTRY are typically much more timely than NCBI
5. NCBI does not cover all of the USPTO granted patent sequence data available in USGENE



## Total sequence count comparison by patent authority in DGENE versus NCBI

31



(Statistics: January 21<sup>st</sup>, 2006)

STN

FIZ Karlsruhe

## In addition, NCBI/EMBL/DDBJ patent records only have minimal bibliographic and text data

32

General Information			
Accession #	AAA00521		
SRS Entry ID	USPO_PRT:AAA00521		
Molecule Type	PRT		
Sequence Length	40		
Entry Data Class	STANDARD		
Sequence Version	AAA00521.1		
Creation Date	21-MAY-1993		
UniParc	<a href="#">UPI0000035113</a>		
Description			
Description	Sequence 1 from Patent US 4563352.		
Organism	Unknown		
References			
1.	Rivier, J.E.F.; Spiess, J. and Vale, W.W. Jr.; <b>Human pancreatic GRF</b> Patent number <a href="#">US4563352-A/1</a> 07-JAN-1986; The Salk Institute For Biological Studies; San Diego, CA Position 1-40		
Features			
Key	Location	Qualifier	Value
<a href="#">source</a>	1..40		
Sequence			
Characteristics	Length: 40 AA		
Sequence	>uspo_prt AAA00521 AAA00521 Sequence 1 from Patent US 4563352. YADAIFTNSYRKVLGQLSARKLLQDIMSRQQGSENGERGA		

**Reminder:** NCBI/EMBL/DDBJ cover sequences from U.S. granted patents – sequences from U.S. published applications are **not** covered (see slide 18).

STN

FIZ Karlsruhe



## Agenda

33

- Sequences in patent publications
- STN patent sequence databases
- Comparison to web based resources
- Results of a comparative search example
- Summary and resources



## A simple BLAST example shows the importance of using STN databases

34

### **Search Question:**

Find patent all references to Breast Cancer 1, early onset isoform 1 (NCBI: NP\_009225), or other very similar proteins (i.e. >80% match)

(Search conducted on June 12<sup>th</sup>, 2008)



# Homo sapiens Breast Cancer 1, early onset isoform 1 protein (NCBI: NP\_009225)

35

The image shows a screenshot of a web browser displaying the NCBI protein viewer for BRCA1 protein (NP\_009225). The protein sequence is shown in a Notepad window, starting with MDLSALRVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKF... The browser interface includes search bars and navigation options.

STN

FIZ Karlsruhe

# The best answer from DGENE...

36

```

=> D BIB AB PSL ORGN SCORE ALIGN
L2 ANSWER 1 OF 114 DGENE COPYRIGHT 2008 THOMSON
AN AFU71404 protein DGENE
TI Diagnosing breast cancer or a susceptibility to breast cancer
   comprises detecting BRCA2 999del5 and BARD1 Cys557Ser.
IN Stacey S N; Sulem P; Thorsteinsdottir U; Kong A
PA (STAC-1) STACEY S N. . . . .
PI US 2007092900 A1 20070426
AI US 2006-515368 20060831
PRAI US 2005-730703P 20051026
OS 2007-556498 [54]
CR N-PSDB: AFU71403
DESC Human breast cancer 1 (BRCA1) protein.
PSL Disclosure; SEQ ID NO 2
AB The present invention relates to a methods for diagnosing and
   characterizing breast cancer or a susceptibility to breast cancer in
   an individual comprising detecting breast cancer 2 (BRCA2) 999del5
   and BRCA1 associated RING domain 1 (BARD1) Cys557Ser. . . . .
ORGN Homo sapiens.
SCORE 3722 100% of query self score 3722
BLASTALIGN
  Query = 1863 letters
  Length = 1863
  Score = 3722 bits (9652), Expect = 0.0
  Identities = 1863/1863 (100%), Positives = 1863/1863 (100%)
  Query: 1 MDLSALRVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKF...
  Sbjct: 1 MDLSALRVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKF...
  
```

114 sequence hits were found in DGENE.

DGENE records feature extensive Thomson value added content.

STN

FIZ Karlsruhe

## The best answer from USGENE...

37

```
=> D BIB AB ECLM ORGN SCORE ALIGN
L4 ANSWER 1 OF 70 USGENE COPYRIGHT 2008 SEQUENCE
AN 20070083334.873670 Protein USGENE
TI Methods and systems for annotating biomolecular sequences
   (PublishedApplication)
IN Mintz Liat (Kendall-Park, NJ); Xie Hangqing (Plainsboro, NJ); . . . .
PA Compugen Ltd
PI US 20070083334 A1 20070412
AI US 2006-443428 20060531
DT Patent
AB A method of annotating biomolecular sequences. The method comprises
   (a) computationally clustering the biomolecular sequences according
   to a progressive homology range, to thereby generate a plurality of
   clusters each being of a predetermined homology . . . .
ECLM US20070083334 A1: 1-98. (canceled) 99. An isolated protein, having a
     sequence according to the amino acid sequence of SEQ ID NO: 836875.
ORGN Homo Sapiens
SCORE 3722 100% of query self score 3722
BLASTALIGN
  Query = 1863 letters
  Length = 1863
  Score = 3722 bits (9652), Expect = 0.0
  Identities = 1863/1863 (100%), Positives = 1863/1863 (100%)
  Query: 1 MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
           MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
  Sbjct: 1 MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
```

70 sequence hits were found in USGENE.

Note that this is a different best answer from that found in DGENE.

STN

FIZ Karlsruhe

## The best answer from PCTGEN...

38

```
=> D BIB ORGN SCORE ALIGN
L6 ANSWER 1 OF 7 PCTGEN COPYRIGHT 2008 WIPO on S
AN 2007047796.8180 PRT PCTGEN
TI TISSUE- AND SERUM-DERIVED GLYCOPROTEINSAND METHODS OF THEIR USE
PA Institute for Systems Biology
   Zhang, Hui
   Aebersold, Rudolf H.
PI WO 2007047796 20070426
RLI US 2005-728044P 20051017
ED 20070427
DT Patent
ORGN Homo sapiens
SCORE 3722 100% of query self score 3722
BLASTALIGN
  Query = 1863 letters
  Length = 1863
  Score = 3722 bits (9652), Expect = 0.0
  Identities = 1863/1863 (100%), Positives = 1863/1863 (100%)
  Query: 1 MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
           MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
  Sbjct: 1 MDLSALRVVEEVQNVINAMQKILECPICLELIKEPVSTKCDHIFCKFCMLKLLNQQKGPSQ
```

7 sequence hits were found in PCTGEN.

Note that this is a different best answer from those found in DGENE and USGENE (previous slides).

STN

FIZ Karlsruhe

## The best patent answer from REGISTRY...

39

Alignment Details

3722 0.0 (934769-03-6) 2: PN: US20070092900 SEQID: 2 unclaimed protein

Length = 1863  
Score = 3722 Expect = 0.0  
Identities = 1863/1863 (100%) Positives = 1863/1863 (100%)

Query: 1 MDLSALRVEEVQNVINAMQKILECPICLELIK  
MDLSALRVEEVQNVINAMQKILECPICLELIK  
Subject: 1 MDLSALRVEEVQNVINAMQKILECPICLELIK

Query: 56 KGPSQCPLCKNDITKRSLOESTRFSQVVEELLKIICAFQLDGTGLEYSYNFYAKK 110  
KGPSQCPLCKNDITKRSLOESTRFSQVVEELLKIICAFQLDGTGLEYSYNFYAKK  
Subject: 56 KGPSQCPLCKNDITKRSLOESTRFSQVVEELLKIICAFQLDGTGLEYSYNFYAKK 110

Query: 111 ENNSPEHLKDEVSIIQSMGYRMRKRLQSEPENPSLQETSLSVQLSNLGTVRTL 165  
ENNSPEHLKDEVSIIQSMGYRMRKRLQSEPENPSLQETSLSVQLSNLGTVRTL  
Subject: 111 ENNSPEHLKDEVSIIQSMGYRMRKRLQSEPENPSLQETSLSVQLSNLGTVRTL 165

Query: 166 RTKQRIQPQKTSVYIELGSDSSEDTVMKATYCSVGDQEL  
RTKQRIQPQKTSVYIELGSDSSEDTVMKATYCSVGDQEL  
Subject: 166 RTKQRIQPQKTSVYIELGSDSSEDTVMKATYCSVGDQEL

Query: 221 AKKAAEFSETDVTNTEHHQPSMNDLNTTEKRAAERHPEKYQSSVSNLHVEPCG 275

Note that this is the same best answer found in DGENE.

53 sequence hits were found in REGISTRY.

STN

FIZ Karlsruhe

## The best patent answer from NCBI...

40

> gb|AAC11849.1|I78815 Sequence 2 from patent US 5693473

gb|AAC20811.1|I81570 Sequence 2 from patent US 5709999

gb|AAC20817.1|I81576 Sequence 2 from patent US 5710001

gb|AAC86103.1|AR005620 Sequence 2 from patent US 574728

gb|AAC87698.1|AR008672 Sequence 2 from patent US 575344

gb|AAE61797.1| Sequence 2 from patent US 6162897

gb|AAT26387.1| Sequence 2 from patent US 6720158

gb|AAV20838.1| Sequence 1098 from patent US 6753314

Length=1863

Score = 3844 bits (9969), Expect = 0  
Identities = 1863/1863 (100%), Positi-

Query 1 MDLSALRVEEVQNVINAMQKILECP  
Sbjct 1 MDLSALRVEEVQNVINAMQKILECPLEDEIKFVSTKCDHIFCKFCMLKLNQKKGPSQ 60

Query 61 CPLCKNDITKRSLOESTRFSQVVEELLKIICAFQLDGTGLEYSYNFYAKKENSPEHLK 120  
Sbjct 61 CPLCKNDITKRS

Query 121 EVSIIQSMGYRN  
Sbjct 121 EVSIIQSMGYRN

Query 181 ELGSDSSEDTVMKATYCSVGDQELDQITFGTRDEISLDSAKKAAEFSETDVTNTEHHQ 240  
Sbjct 181 ELGSDSSEDTVMKATYCSVGDQELDQITFGTRDEISLDSAKKAAEFSETDVTNTEHHQ 240

Query 241 PSNNDLNTTEKRAAERHPEKYQSSVSNLHVEPCGNTTHASSLQHENSSLLLTKDRMNVE 300  
Sbjct 241 PSNNDLNTTEKRAAERHPEKYQSSVSNLHVEPCGNTTHASSLQHENSSLLLTKDRMNVE 300

32 sequence hits were found at NCBI.

Note that these patents are not the same best answers seen in DGENE, USGENE, PCTGEN or REGISTRY.

Reminder: NCBI only provides protein patent sequence coverage from U.S. issued patents. U.S. and WIPO/PCT (WO) published application protein sequences cannot be retrieved at the NCBI.

STN

FIZ Karlsruhe

## Summary of results for Breast Cancer 1, early onset isoform 1 protein (NCBI: NP\_009225)

41

	<b>SEQs &gt; 80%</b>	<b>PNs</b>	<b>Patent Families</b>
<b>DGENE</b>	114	36	26
<b>USGENE</b>	70	38	22
<b>PCTGEN</b>	7	4	4
<b>REGISTRY</b>	52	30	23
<b>NCBI*</b>	32	21	9
<b>TOTAL</b>	-	-	34

(\* All of the NCBI patent hits were for U.S. patents, and none of them were unique compared to USGENE.)



## Agenda

42

- Sequences in patent publications
- STN patent sequence databases
- Comparison to web based resources
- Results of a comparative search example
- Summary and resources



## Comparing databases...

43


- DGENE
  - The most comprehensive global patent sequence database
  - Implemented in-house at major patent offices\*
- USGENE
  - Sequences from equivalent USPTO applications and patents
  - Incorporates all U.S. patent sequence data available at NCBI
- PCTGEN
  - Sequences from equivalent WIPO/PCT publications
  - The complete collection of e-published sequences from WIPO
- REGISTRY
  - Complementary value-added indexing to DGENE
  - Unique non-patent literature coverage

\* See page 11: [http://www.trilateral.net/projects/biotechnology/search\\_guidebook\\_vers\\_1.pdf](http://www.trilateral.net/projects/biotechnology/search_guidebook_vers_1.pdf)



## Database timeliness is an important factor in understanding comprehensiveness

44

	Update Frequency	Typical Timeliness	Value added
PCTGEN	Weekly	24 hours	
USGENE	Weekly	3 days	
REGISTRY	Daily	27 days	
DGENE	Biweekly	65 days	
NCBI/EMBL	Daily	1-6 months	



## Resources for sequence searching on STN 45

- *Sequence Searching on STN* modular workshop  
<http://www.fiz-k.com/bostonsequenceworkshop>
  - Sequence Code Match (SCM) searching
  - DGENE, USGENE, PCTGEN content and searching
  - CAS REGISTRY and REGISTRY BLAST
  - Multifile searching using USGENE and DGENE
- USGENE resources, reference materials and FAQ  
<http://www.sequencebase.com>
- CAS REGISTRY sequence coverage and resources  
<http://www.cas.org/support/stngen/stndoc/sequences.html>

**STN**

 FIZ Karlsruhe

**STN<sup>®</sup>**

Resources for Intellectual Property  
Sequence searching on STN

[www.stn-international.com](http://www.stn-international.com)