

# Project Description

Laboratory of Bioinformatics I  
Module 2

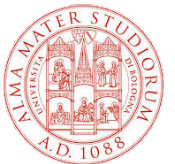
Emidio Capriotti

<http://biofold.org/>



**Biomolecules**  
**Folding and**  
**Disease**

Department of Pharmacy  
and Biotechnology (FaBiT)  
University of Bologna



# Main Aim

Building a Profile Hidden Markov Model for the Kunitz-type protease inhibitor domain.

Kunitz domains are the active domains of proteins that **inhibit the function of protein degrading enzymes** or, more specifically, domains of Kunitz-type are **protease inhibitors**.

Examples of Kunitz-type protease inhibitors are aprotinin (bovine pancreatic trypsin inhibitor, BPTI), Alzheimer's amyloid precursor protein (APP), and tissue factor pathway inhibitor (TFPI).

# Aprotinin

The drug **aprotinin** (Trasylol, previously Bayer and now Nordic Group pharmaceuticals), is the small protein **bovine pancreatic trypsin inhibitor** (BPTI), an **antifibrinolytic** molecule that inhibits trypsin and related proteolytic enzymes. Under the trade name Trasylol, aprotinin was used as a medication administered by injection **to reduce bleeding** during complex surgery, such as heart and liver surgery. Its main effect is the slowing down of fibrinolysis, the process that leads to the breakdown of blood clots. The aim in its use was to decrease the need for blood transfusions during surgery, as well as end-organ damage due to hypotension (low blood pressure) as a result of marked blood loss.

BPTI is the classic member of the protein family of Kunitz-type serine protease inhibitors. Its physiological functions include the **protective inhibition of the major digestive enzyme trypsin** when small amounts are produced by cleavage of the trypsinogen precursor during storage in the pancreas.

# Aprotinin Structure

Aprotinin is a **monomeric** (single-chain) globular polypeptide derived from bovine lung tissue. It has a molecular weight of 6512 and consists of a chain 58 residues long that folds into a **stable, compact tertiary structure of the 'small SS-rich' type, containing 3 disulfides, a twisted  $\beta$ -hairpin and a C-terminal  $\alpha$ -helix.**

There are 10 positively-charged lysine (K) and arginine (R) side chains and only 4 negative aspartate (D) and glutamates (E), making the protein strongly basic

The high stability of the molecule is due to the **3 disulfide bonds linking the 6 cysteine members of the chain (Cys5-Cys55, Cys14-Cys38 and Cys30-Cys51).**

**The long, basic lysine 15 side chain on the exposed loop binds very tightly in the specificity pocket at the active site of trypsin and inhibits its enzymatic action.** BPTI is synthesized as a longer, precursor sequence, which folds up and then is cleaved into the mature sequence.

# Start from the Structure

In the Protein Data Bank the crystal of 3TGI a complexed of the BPTI

Structure Summary | 3D View | Annotations | Experiment | Sequence | Genome | Versions

Biological Assembly 1 ?

Display Files | Download Files

## 3TGI

WILD-TYPE RAT ANIONIC TRYPSIN COMPLEXED WITH BOVINE PANCREATIC TRYPSIN INHIBITOR (BPTI)

DOI: [10.2210/pdb3TGI/pdb](https://doi.org/10.2210/pdb3TGI/pdb)

Classification: **COMPLEX (SERINE PROTEASE/INHIBITOR)**

Organism(s): [Rattus norvegicus](#), [Bos taurus](#)

Mutation(s): No ⓘ

Deposited: 1998-07-15 Released: 1998-12-23

Deposition Author(s): [Pasternak, A.](#), [Ringe, D.](#), [Hedstrom, L.](#)

**Experimental Data Snapshot**

Method: X-RAY DIFFRACTION

Resolution: 1.80 Å

R-Value Free: 0.210

R-Value Work: 0.177

R-Value Observed: 0.177

**wwPDB Validation** ⓘ [3D Report](#) [Full Report](#)

Metric	Percentile Ranks	Value
Rfree		0.193
Clashscore		4
Ramachandran outliers		0.4%
Sidechain outliers		4.3%
RSRZ outliers		0.7%

Worse | Better

■ Percentile relative to all X-ray structures  
□ Percentile relative to X-ray structures of similar resolution

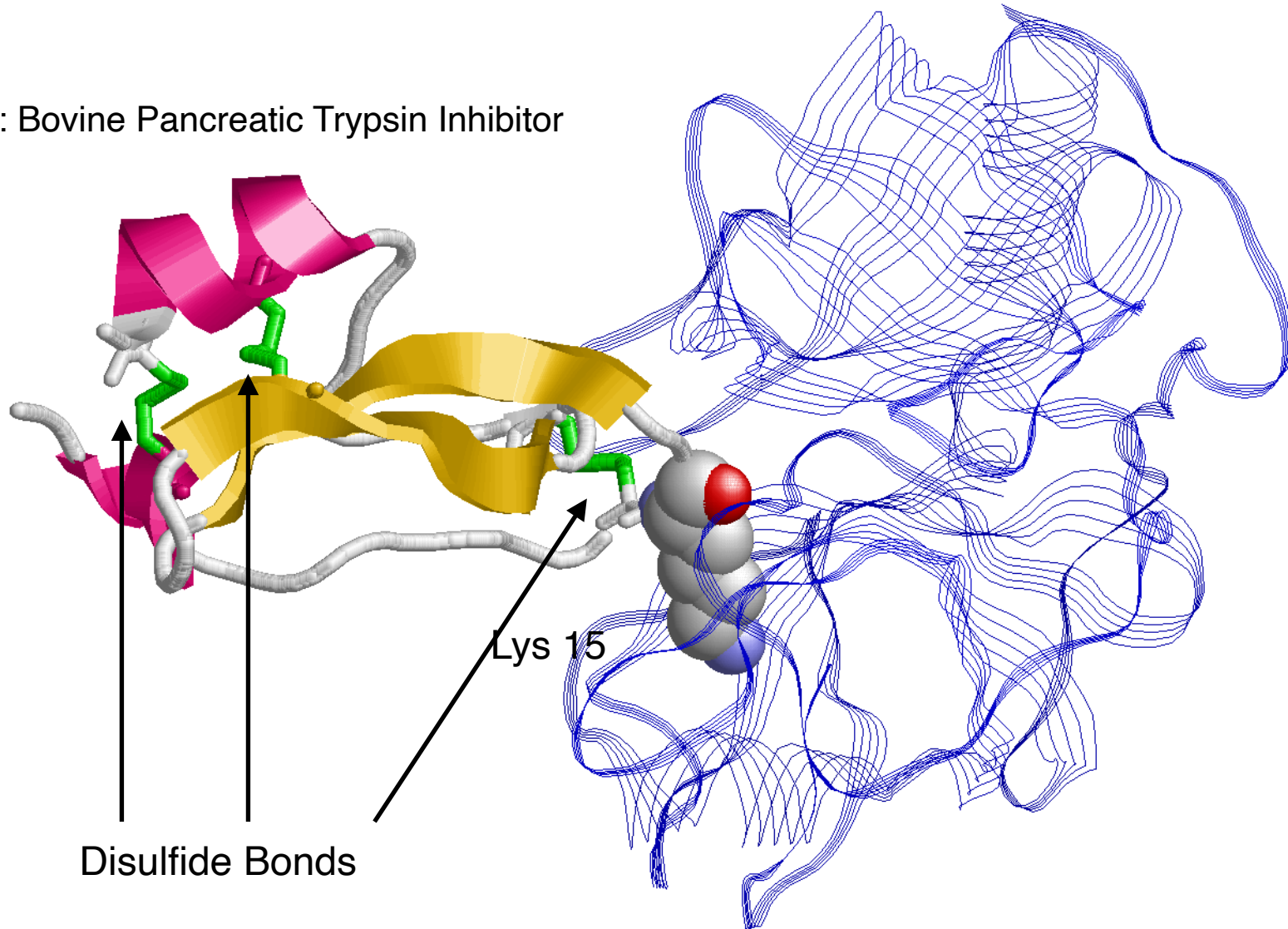
Find Similar Assemblies

This is version 1.2 of the entry. See complete [history](#).

# Structure Analysis

Chain E: Rat Anionic Trypsin

Chain I: Bovine Pancreatic Trypsin Inhibitor



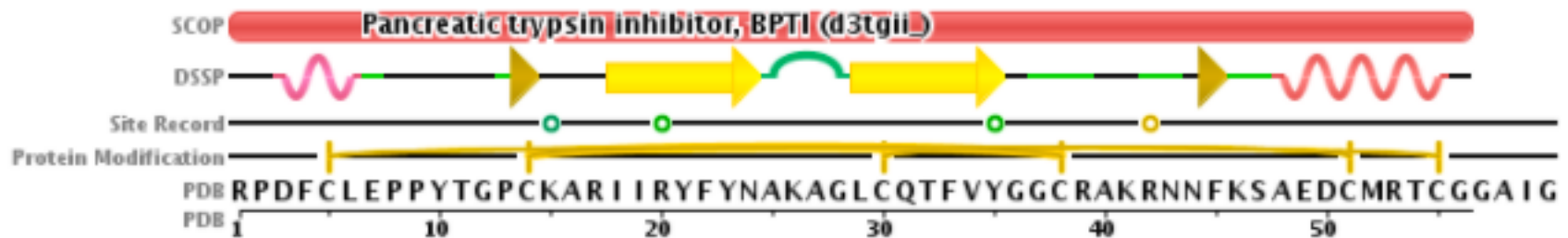
# The Protein Fold

The **structure is a disulfide rich alpha+beta fold**. Bovine pancreatic trypsin inhibitor is an extensively studied model structure.

The majority are restricted to metazoa with a single exception: *Amsacta moorei entomopoxvirus*, a species of poxvirus.

**They are short (about 50 to 60 amino acid residues) alpha/beta proteins with few secondary structures. The fold is constrained by three disulfide bonds.**

## Sequence Chain View



# Annotation

In UniProt we found the information about the function and important sites

**UniProt** BLAST Align Peptide search ID mapping SPARQL UniProtKB  Advanced | List Search Help

**Function** **P00974 · BPT1\_BOVIN**

**Protein<sup>i</sup>** Pancreatic trypsin inhibitor **Amino acids** 100

**Status<sup>i</sup>** UniProtKB reviewed (Swiss-Prot) **Protein existence<sup>i</sup>** Evidence at protein level

**Organism<sup>i</sup>** **Bos taurus (Bovine)** **Annotation score<sup>i</sup>** 5/5

[Entry](#) [Feature viewer](#) [Publications](#) [External links](#) [History](#)

BLAST [Download](#) [Add](#) [Add a publication](#) [Entry feedback](#)

**Function<sup>i</sup>**  
Inhibits trypsin, kallikrein, chymotrypsin, and plasmin.

**Features**  
Showing features for site<sup>i</sup>.

1 10 20 30 40 50 60 70 80 90 100

MKMSRLCLSLVALLVLLGLTAASTPGCDTSNQAQAKAQRPDFCLEPPYTGPKARIIRYFYNAKAGLCQTFVYGGCRKRNNPKSAEDCMRTC GGAIGPWENL

TYPE	ID	POSITION(S)	DESCRIPTION
-- Select --			
▶ Site		50-51	Reactive bond for trypsin


**GO annotations<sup>i</sup>** [Expand table](#)

[Feedback](#) [Help](#)



# PFAM

The Kunitz BPTI family is described in PFAM database



InterPro - Member

Classification of protein families

Home Search Browse Results Release notes Download Help About Contact us

Home / Browse / By Entry / Pfam / PF00014 / Overview

## Pfam PF00014 Kunitz/Bovine pancreatic trypsin inhibitor domain

Pfam entry ⓘ

Overview	
Proteins	31k
Domain Architectures	3k
Taxonomy	5k
Proteomes	1k
Structures	210
Signature	
AlphaFold	20k
Alignment	
Curation	

Member database [Pfam](#) ⓘ

Pfam type domain

Short name *Kunitz\_BPTI*

### Description

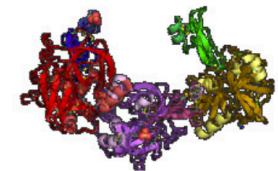
Indicative of a protease inhibitor, usually a serine protease inhibitor. Structure is a disulfide rich alpha+beta fold. BPTI (bovine pancreatic trypsin inhibitor) is an extensively studied model structure. Certain family members are similar to the tick anticoagulant peptide (TAP, Swiss:P17726). This is a highly selective inhibitor of factor Xa in the blood coagulation pathways [2]. TAP molecules are highly dipolar [1], and are arranged to form a twisted two- stranded antiparallel beta-sheet followed by an alpha helix [2].

[Add your annotation](#) ▾

Integrated to

[> IPR002223](#)

Representative structure



# Domain Organization

The Kunitz domain is present in many proteins with different architectures

**InterPro** - Member  
Classification of protein families

Home | Search | **Browse** | Results | Release notes | Download | Help | About | Contact us

Home / Browse / By Entry / Pfam / PF00014 / Domain Architecture

## Pfam PF00014 Kunitz/Bovine pancreatic trypsin inhibitor domain

[Pfam entry](#)

**2677 domain architectures**

Overview  
Proteins 31k  
**Domain Architectures 3k**  
Taxonomy 5k  
Proteomes 1k  
Structures 210  
Signature  
AlphaFold 20k  
Alignment  
Curation

[Export](#)

There are 8227 proteins with this architecture (represented by P00975):  
PF00014  
● Kunitz\_BPTI  
60

There are 2911 proteins with this architecture (represented by Q20014):  
PF00014 - PF00014  
● ● Kunitz\_BPTI  
219

There are 2055 proteins with this architecture (represented by P10646):  
PF00014 - PF00014 - PF00014  
● ● ● Kunitz\_BPTI  
304

# PFAM Alignments

PFAM stores different alignments with increasing number of sequences.

**InterPro** - Member  
Classification of protein families

Home | Search | **Browse** | Results | Release notes | Download | Help | About | Contact us

Home / Browse / By Entry / Pfam / PF00014 / Entry Alignments

## **Pfam** PF00014 Kunitz/Bovine pancreatic trypsin inhibitor domain

Pfam entry ⓘ

Available alignments: seed (99)

Colors: clustal2 Conservation:  **Legends** **Download**

**99 sequences**

Overview	Proteins 31k
Domain Architectures 3k	Taxonomy 5k
Proteomes 1k	Structures 210
Signature	AlphaFold 20k
<b>Alignment</b>	Curation

017644\_CAEEL/982-1034  
A8XY36\_CAEER/1446-1497  
Q94164\_CAEEL/22-79  
Q18761\_CAEEL/35-87  
EPPI\_HUMAN/76-128  
O45881\_CAEEL/1082-1134  
AMBP\_BOVIN/286-338  
Q21418\_CAEEL/410-462  
Q23456\_CAEEL/256-309  
PPN1\_CAEEL/1852-1904  
PPN1\_CAEEL/1913-1965  
CO7A1\_HUMAN/2875-2930  
PPN1\_CAEEL/1270-1322  
O45881\_CAEEL/1222-1274  
CO6A3\_HUMAN/3111-3163

```
FCL SAR . DSG . P . CN . . . N . . . . . FE . . KRYGYD . . . . . ANTDTCVEYQYGGCEGT . . . L . . . . . NNFHSLQRCTEIK
VCDEAK . DTG . P . CT . . . N . . . . . FA . . . TKWYYN . . . . . QADGTCNRFHYGGCQGT . . . N . . . . . NFRDNEQQCKAAK
RCSKSI FDSNLTAKCE . . . . . KSS . TIKFHFD . . . . . QSTGLCMNFRWDGCKDQ . . . E . . . . . NKFDLSQECASG
ICLEDV . DPG . P . CQ . . . Y . . . . . YQ . . . VQWFD . . . . . KQVEECKEFHYGGCMGT . . . K . . . . . NRFSSKQQCVKQK
VCEMPK . ETG . P . CL . . . A . . . . . YF . . . LHWWYD . . . . . KKDNTCSMFVYGGCQGN . . . N . . . . . NNFQSKANCLNTG
KCLQPV . EPG . P . CK . . . N . . . . . FA . . . DRWYFN . . . . . VDDGTCHPFKYGGCAGN . . . R . . . . . NHFFTQKECEVHG
ACNLP I . VQG . P . CR . . . S . . . . . YI . . . QLWAFD . . . . . AVKGKCVRFYGGCCKGN . . . G . . . . . NKFYSEKECKEYK
VCKLPR . EQG . N . CG . . . T . . . . . YS . . . NRWFFN . . . . . AKTGNCEEF IYSGCQGN . . . A . . . . . NNFETYKECQDYK
PCSLSP . DKG . FP . GS . . . V . . . . . TV . . . NMWYYD . . . . . PTSTTCSPFMYLKGKGN . . . S . . . . . NRFETSEECLETG
FCTLER . SAG . P . CT . . . D . . . . . SI . . . SMWYFD . . . . . STHLDCKPFTYGGCRGN . . . Q . . . . . NRFVSKEQCQSSK
ICTLRP . EPG . P . CR . . . L . . . . . GL . . . EKIFYD . . . . . PVIQSCHMFHYGGCEGN . . . A . . . . . NRFDSSELD CFRRK
PCSLPL . DEG . S . CT . . . A . . . . . YT . . . LRWYHRA . . . . . VTGSTEACHPFVYGGCGGN . . . L . . . . . NRFGTREACERRK
ICRSRQ . DAG . P . CE . . . T . . . . . YS . . . DQWFYN . . . . . AFSQECETFTYGGCGGN . . . A . . . . . NRRFRSKDECEQRK
VCAMPP . DAG . V . CT . . . N . . . . . YT . . . PRWFFN . . . . . SQTQQCEQFAYGSCGGN . . . E . . . . . NNFDRNTCERKIK
ICKLPK . DEG . T . CR . . . D . . . . . FI . . . LKWYYD . . . . . PNTKSCARFWYGGCGGN . . . E . . . . . NKFGSQKECEKVI
```

# PFAM Curation

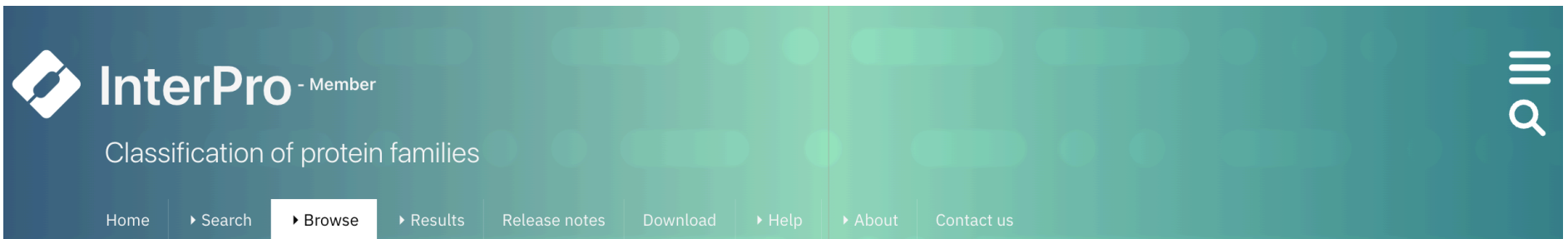
Information about the PFAM family alignment is reported in the Curation page

The screenshot shows the InterPro website interface. At the top, the InterPro logo and tagline 'Classification of protein families' are visible. A navigation menu includes 'Home', 'Search', 'Browse', 'Results', 'Release notes', 'Download', 'Help', 'About', and 'Contact us'. The breadcrumb trail reads: Home / Browse / By Entry / Pfam / PF00014 / Curation. The main heading is 'Pfam PF00014 Kunitz/Bovine pancreatic trypsin inhibitor domain', with a 'Pfam entry' link. A left sidebar contains a navigation menu with categories: Overview, Proteins (31k), Domain Architectures (3k), Taxonomy (5k), Proteomes (1k), Structures (210), Signature, AlphaFold (20k), Alignment, and Curation (selected). The main content area is titled 'Curation' and contains the following information:

- Author:** Fenech M
- Sequence Ontology:** SO:0000417
- HMM Information:**
  - HMM build commands:** Build method: hmmbuild -o /dev/null HMM SEED; Search method: hmmsearch -Z 75585367 -E 1000 --cpu 4 HMM pfamseq
  - Gathering threshold:** Sequence: 21; Domain: 21
- Download:** Download the raw HMM for this family

# HMM Logo

Important protein sites can be visualized using HMM Logo



InterPro - Member  
Classification of protein families

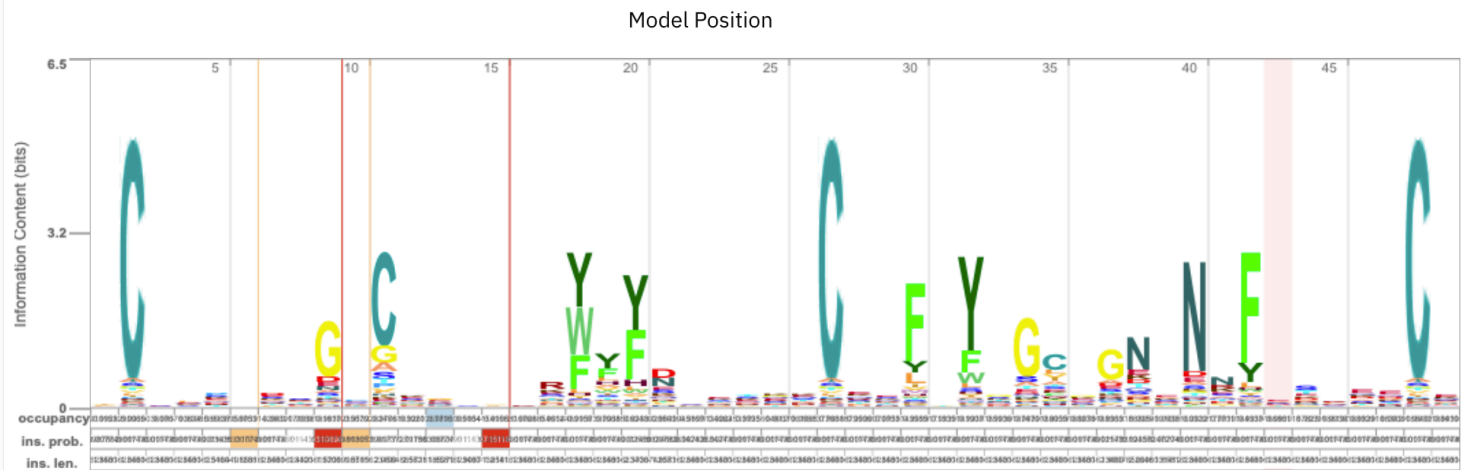
Home | Search | Browse | Results | Release notes | Download | Help | About | Contact us

Home / Browse / By Entry / Pfam / PF00014 / Logo

## Pfam PF00014 Kunitz/Bovine pancreatic trypsin inhibitor domain

[Pfam entry](#)

- Overview
- Proteins 31k
- Domain Architectures 3k
- Taxonomy 5k
- Proteomes 1k
- Structures 210
- Signature**
- AlphaFold 20k
- Alignment
- Curation



# Specific Aims

The specific aims are:

1. Build your own model for the Kunitz domain, starting from available structural information.
2. Use the model for annotating Kunitz domains in SwissProt.

Write a detailed draft of the project identifying

- the main steps;
- the sources of the data to be analyzed;
- the procedures/programs you would adopt;
- The results to be produced for validating your model

# Structure Selection

## Retrieve available structures of the Kunitz domain

This is the crucial step: you need to collect a large set of structures that are endowed with

**Source:** *PDB*

**Method:** *different alternative options are possible:*

- a. consider a prototype structure and search in the PDB other similar structures (e.g, by using the PDBe-fold web site)*
- b. retrieve from UniProt the protein endowed with an annotated BPTI/ Kunitz type domain and with a 3D structure covering it.*
- c. Try to directly scan the PDB for structurally-resolved Kunitz domains(e.g., you can use the CATH code 4.10.410.10)*
- d. ....*

# Possible Issues

When **selecting the domains** for building the seed alignment, keep in mind that:

- PDB files can contain **more than one chain**;
- **A chain can contain different domains** of the same type or of different types;
- Structures of the same protein can be found in **different PDB files**;
- the PDB collects the structure of **mutated proteins**;
- **Resolution** can be an issue during structural alignment.



# Protein Alignment

Perform the structural alignment of the selected domains

**Method:** Any multiple structural alignment method (e.g. PDBe-fold)

*On the basis of the structural alignment results you can correct/refine your initial choice of the seed proteins.*

If needed convert the alignment in Stockholm format

**Method:** JalView or write an ad-hoc program

# Generate HMM Model

Train a profile HMM

**Method:** *HMMER hmmbuild routine*

Verify that the trained HMM is able to recognize the proteins in your dataset  
(consistency test)

**Method:** *HMMER hmmsearch routine*

If the performance on the train set is low there is probably some problem in the set of proteins your choose and/or in the alignment you fed to HMM during the training procedure

# Method Testing

## Retrieve a suitable dataset for validating the HMM prediction

Only manually curated proteins should be considered, avoiding fragments  
The dataset should be divided into proteins containing or not containing the BPTI/Kunitz domain (the positive test set should exclude the training data).

*Source: UniProt/Swiss-Prot*

**Method:** *The “advanced search” interface in UniProt web site*

*Different “Gold standard” for defining the positive class are possible:*

- a) the presence of an annotated BPTI/Kunitz domain in the Uniprot entry*
- b) the presence of an annotated PF00014 PFAM domain*
- c) ..*

## Search the validation dataset against the trained model

**Method:** *HMMER hmmsearch routine*

## Compute the scoring indexes for evaluating your profile HMM on the validation sets

**Method:** *Write a program that compares the prediction with the “real” annotations, computes a confusion matrix and the scoring indexes.*

# Analyze the Results

Analyze the results and try to understand whether it is possible to improve them

Prediction could be in some cases optimized by changing the E-value threshold or by refining the training alignment.

Discuss the False Positive and the False Negative predictions

Find your domain in all the SwissProt sequences, comment with respect to the available annotations and comment about the distribution of the Kunitz domain

# Project Report

Project description in the “Bioinformatics” style paper

[http://www.oxfordjournals.org/our\\_journals/bioinformatics/for\\_authors/submission\\_online.html](http://www.oxfordjournals.org/our_journals/bioinformatics/for_authors/submission_online.html)

## Structured Abstract (see recent issues of journal for examples)

### *Original papers*

Abstracts are structured with a standard layout such that the text is divided into sub-sections under the following five headings: **Motivation**, **Results**, [Availability and Implementation], **Contact** [and Supplementary Information]. In cases where authors feel the headings inappropriate, some flexibility is allowed. The abstracts should be succinct and contain only material relevant to the headings. **A maximum of 150 words is recommended.**

- *Motivation*: This section should specifically state the scientific question within the context of the field of study.
- *Results*: This section should summarize the scientific advance or novel results of the study, and its impact on computational biology.

# Main Report

## Introduction

The section must describe the problem treated in the paper, the available knowledge on it. Only information relevant within the scope of the paper should be reported. Appropriate references must cited.

## Materials and Methods

The section must contain the description of the adopted dataset and of the methods that have been used and/or implemented, including the validation procedures and the adopted scoring indexes. Adopted choice must be justified. In principle, **it must contain all the information necessary to integrally reproduce the work.**

## Results (and discussion)

The section must present the obtained results, the possible refinements, and the analysis of the strength and the weakness of the method. Discussion (can be a separate section) must report the considerations that can be derived from results, also in relation to the adopted procedures and/or datasets. ....

## Conclusions

The section present concisely the achievements of the presented work.

# Reference and Data

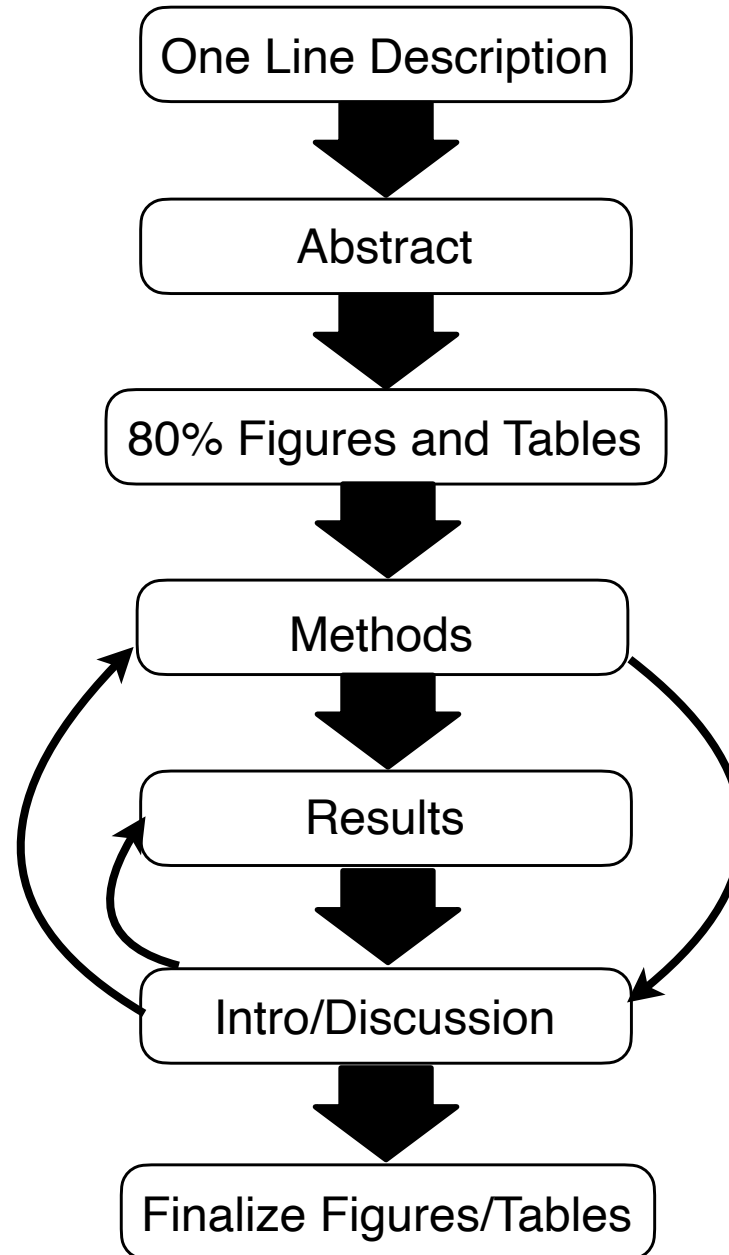
## **References**

See the template for the appropriate format

## **Supplementary Materials**

Supplementary file useful for the presentation of the work can be provided

# Flow Chart





# Project Submission

- The presentation and the approval of the project paper is necessary but not sufficient condition to pass the exam.
- Submit the paper through the following link: <https://bit.ly/biofold-projects>

# Exercise

Build a *blast*-based method to predict the presence of BPTI/Kunitz domain in proteins available in SwissProt using the human proteins as a reference.

- Select all Proteins in SwissProt with BPTI/Kunitz domain.
- Separate human from non human proteins. Use the **non human proteins as a positive** in the testing set.
- Generate a **random set of negative** of the same size of the positive set.
- Remove both positives and negatives from SwissProt and perform the **prediction based on the results of the *blast* search**.