

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 β -hairpin and a C-terminal α -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 given above.

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 ?



3D View: [Structure](#) | [Electron Density](#) | [Ligand Interaction](#)

Global Symmetry: Asymmetric - C1 ?

Global Stoichiometry: Hetero 2-mer - A1B1 ?

[Find Similar Assemblies](#)

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

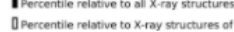
R-Value Work: 0.177

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

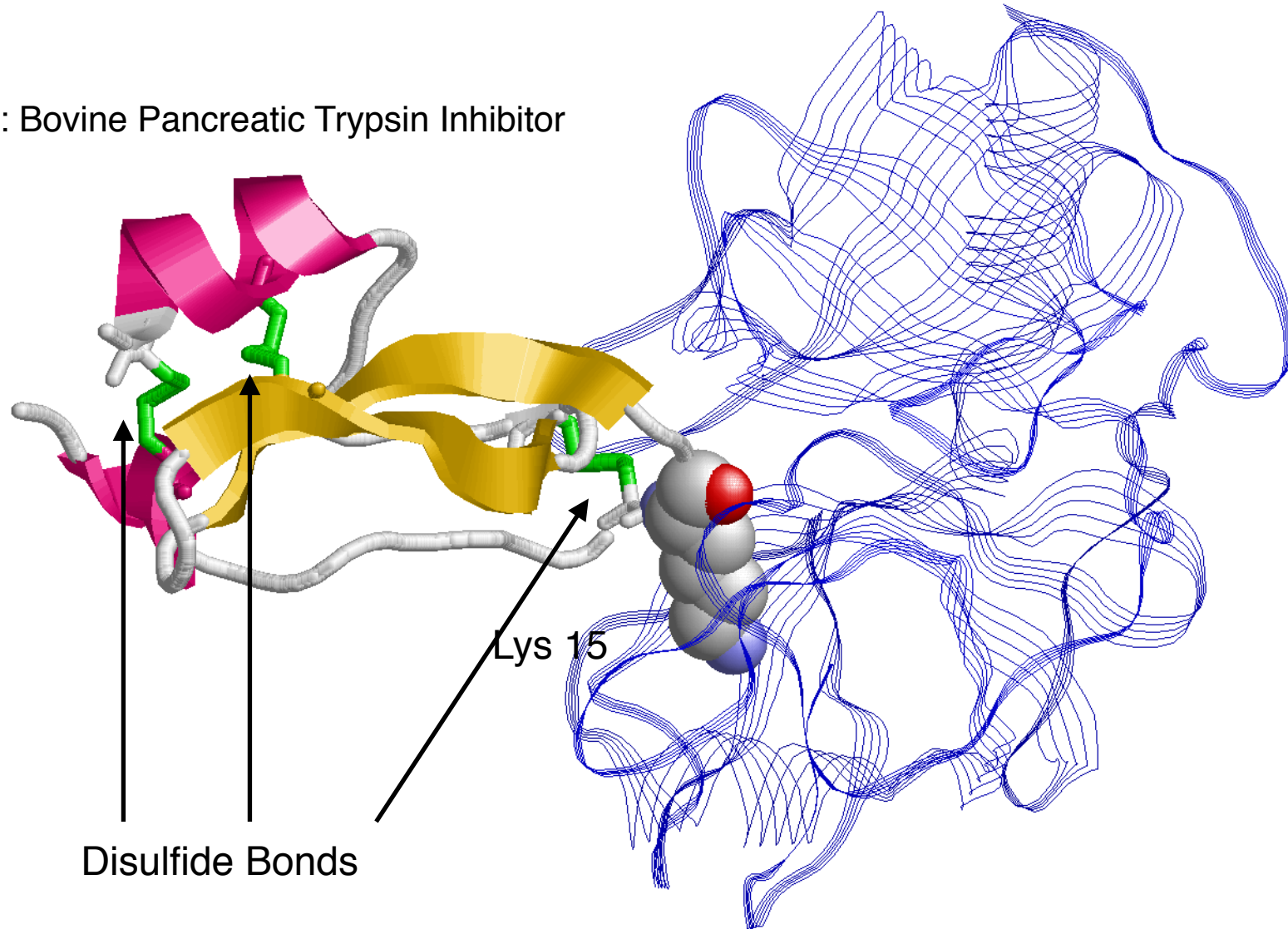
■ Percentile relative to all X-ray structures
□ Percentile relative to X-ray structures of similar resolution

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

Structure Analysis

Chain E: Rat Anionic Trypsin

Chain I: Bovine Pancreatic Trypsin Inhibitor



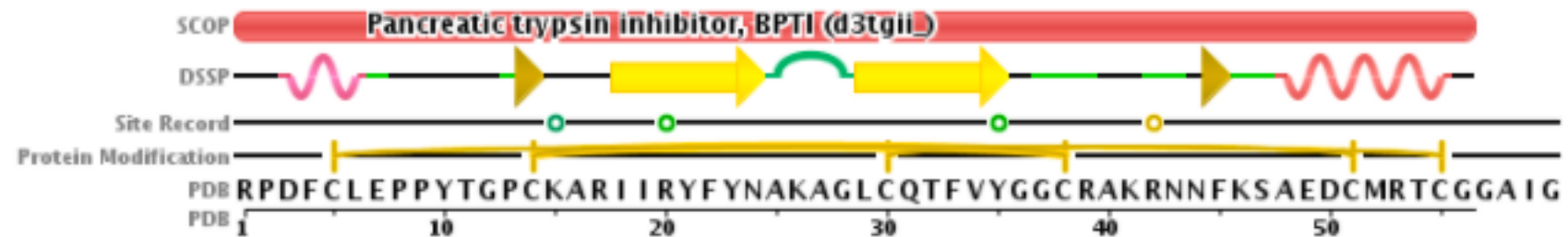
The Protein Fold

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

UniProtKB - P00974 (BPT1_BOVIN)

Basket

Display

Entry

Publications

Feature viewer

Feature table

None

☒ Function

☒ Names & Taxonomy

☒ Subcellular location

☒ Pathology & Biotech

☒ PTM / Processing

☐ Expression

BLAST Align Format Add to basket History

Feedback Help video Other tutorials and videos

Protein | **Pancreatic trypsin inhibitor**

Gene | N/A

Organism | *Bos taurus* (Bovine)

Status | Reviewed - Annotation score: ●●●●● - Experimental evidence at protein levelⁱ

Functionⁱ

Inhibits trypsin, kallikrein, chymotrypsin, and plasmin.

Sites

Feature key	Position(s)	Description	Actions	Graphical view	Length
Site ⁱ	50 – 51	Reactive bond for trypsin			2

PTM / Processingⁱ

Molecule processing


Feature key	Position(s)	Description	Actions	Graphical view	Length
Signal peptide ⁱ	1 – 21	Sequence analysis	Add BLAST		21
Propeptide ⁱ (PRO_0000016852)	22 – 35		Add BLAST		14
Chain ⁱ (PRO_0000016853)	36 – 93	Pancreatic trypsin inhibitor	Add BLAST		58
Propeptide ⁱ (PRO_0000016854)	94 – 100				7

Amino acid modifications

Feature key	Position(s)	Description	Actions	Graphical view	Length
Disulfide bond ⁱ	40 ↔ 90				
Disulfide bond ⁱ	49 ↔ 73	PROSITE-ProRule annotation 1 Publication			
Disulfide bond ⁱ	65 ↔ 86	PROSITE-ProRule annotation 1 Publication			


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](#)

[Home](#) / [Browse](#) / [By Entry](#) / [Pfam](#) / [PF00014](#) / [Overview](#)

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

[Pfam entry](#)

Overview

[Proteins](#) 30k

[Domain Architectures](#) 2k

[Taxonomy](#) 5k

[Proteomes](#) 987

[Structures](#) 194

[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 [PMID:7925983]. TAP molecules are highly dipolar [PMID:10716178], and are arranged to form a twisted two- stranded antiparallel beta-sheet followed by an alpha helix [PMID:7925983].

[Add your annotation](#)

Integrated to
[IPR002223](#)


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](#)

[Home](#) / [Browse](#) / [By Entry](#) / [Pfam](#) / [PF00014](#) / [Domain Architecture](#)

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

[Pfam entry](#)

Overview

Proteins 30k

Domain Architectures 2k

Taxonomy 5k

Proteomes 987

Structures 194


Signature

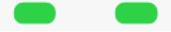
AlphaFold 20k


Alignment

Curation

2456 domain architectures found.

There are 7987 proteins with this architecture (represented by [P00975](#)):
[PF00014](#)
 Kunitz_BPTI
60

There are 2723 proteins with this architecture (represented by [Q20014](#)):
[PF00014](#) - [PF00014](#)
 Kunitz_BPTI
219

There are 1983 proteins with this architecture (represented by [P10646](#)):
[PF00014](#) - [PF00014](#) - [PF00014](#)
 Kunitz_BPTI
304

Show 20 results with this architecture (represented by [P05067](#)):
[PF02177](#) - [PF12924](#) - [PF00014](#) - [PF12925](#) - [PF03494](#) - [PF10515](#)

[Previous](#) [Next](#)

PFAM Alignments

PFAM stores different alignments with increasing number of sequences.

Member

Classification of protein families

Home

Search

Browse

Results

Release notes

Download

Help

About

/ [Browse](#) / [By Entry](#) / [Pfam](#) / [PF00014](#) / [Entry Alignments](#)

Pfam

Kunitz/Bovine pancreatic trypsin inhibitor domain

PF00014

Pfam entry

1

Overview

Proteins 30k

Domain Architectures 2k

Taxonomy 5k

Proteomes 987

Structures 194

Signature

AlphaFold 20k

Alignment

Curation

Available alignments:

seed (99)

Colors: clustal2

Conservation: ☐

Legends

Download

99 Sequences

1

10

20

30

40

50

60

70

80

017644_CAEEL/982-1034

A8XY36_CAEEL/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

FCL

SAR

DSG

P

CN

N

FE

KRYGYD

ANTDT

CVEYQYGGCEGT

L

NNFHS

LQRCTETC

VCDEAK

DTG

P

CT

N

FA

TKWYYN

QADGT

CNRFHYGGCQGT

N

NRFDNE

QQCKAAC

RCSKSI

FDSNLTAKCE

KSS

T

IKFHFD

QSTGL

CMNFRWDGCKDQ

E

NKFDS

LQECASTC

I

CLEDV

DPG

P

CQ

Y

YQ

VQWFD

KQVEECK

EFHYGGCMGT

K

NRFSS

KQQCVKQC

VCEMPK

ETG

P

CL

A

YF

LHWYD

KKDNT

CSMFVYGGCQGN

N

NNFQS

KANCLNTC

KCLQPV

EPG

P

CK

N

FA

DRWYFN

VDDGT

CHPFKYGGCAGN

R

NHFFT

QKECEVHC

ACNLPI

VQG

P

CR

S

YI

QLWAFD

AVK

GKCVRF

SYGGCKGN

G

NKFY

SEKECKEY

C

VCKLPR

EQG

N

CG

T

YS

NRWWFN

AKTGN

CEEFI

YSGCQGN

A

NNFET

YKECQDYC

PCSLSP

DKG

FP

GS

V

TV

NMWYYD

PTSTT

CSPFMYLGKGN

S

NRFET

SEECELET

C

FCTLR

SAG

P

CT

D

SI

SMWYFD

STHLD

CKPFTYGGCRGN

Q

NRFVS

KEQCCQSC

I

CTLRP

EPG

P

CR

L

GL

EKYFYD

PVI

QSCHMFHYGGCEGN

A

NRFDS

ELDCFRRC

PCSLPL

DEG

S

CT

A

YT

LRWYHRA

VTG

STEACH

PFVYGGCGGN

A

NRFGT

REACERRC

I

CRSRQ

DAG

P

CE

T

YS

DQWFYN

AFSQEC

ETFTYGGCGGN


L

NRFRS



KDCECEQRC

PFAM Curation

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


 **InterPro** - Member

Classification of protein families

[Home](#) | [Search](#) | **[Browse](#)** | [Results](#) | [Release notes](#) | [Download](#) | [Help](#) | [About](#)

[Home](#) / [Browse](#) / [By Entry](#) / [Pfam](#) / [PF00014](#) / [Curation](#)

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

[Pfam entry](#) ¹

Overview

Proteins 30k

Domain Architectures 2k

Taxonomy 5k

Proteomes 987

Structures 194

Signature

AlphaFold 20k

Alignment

Curation

Curation

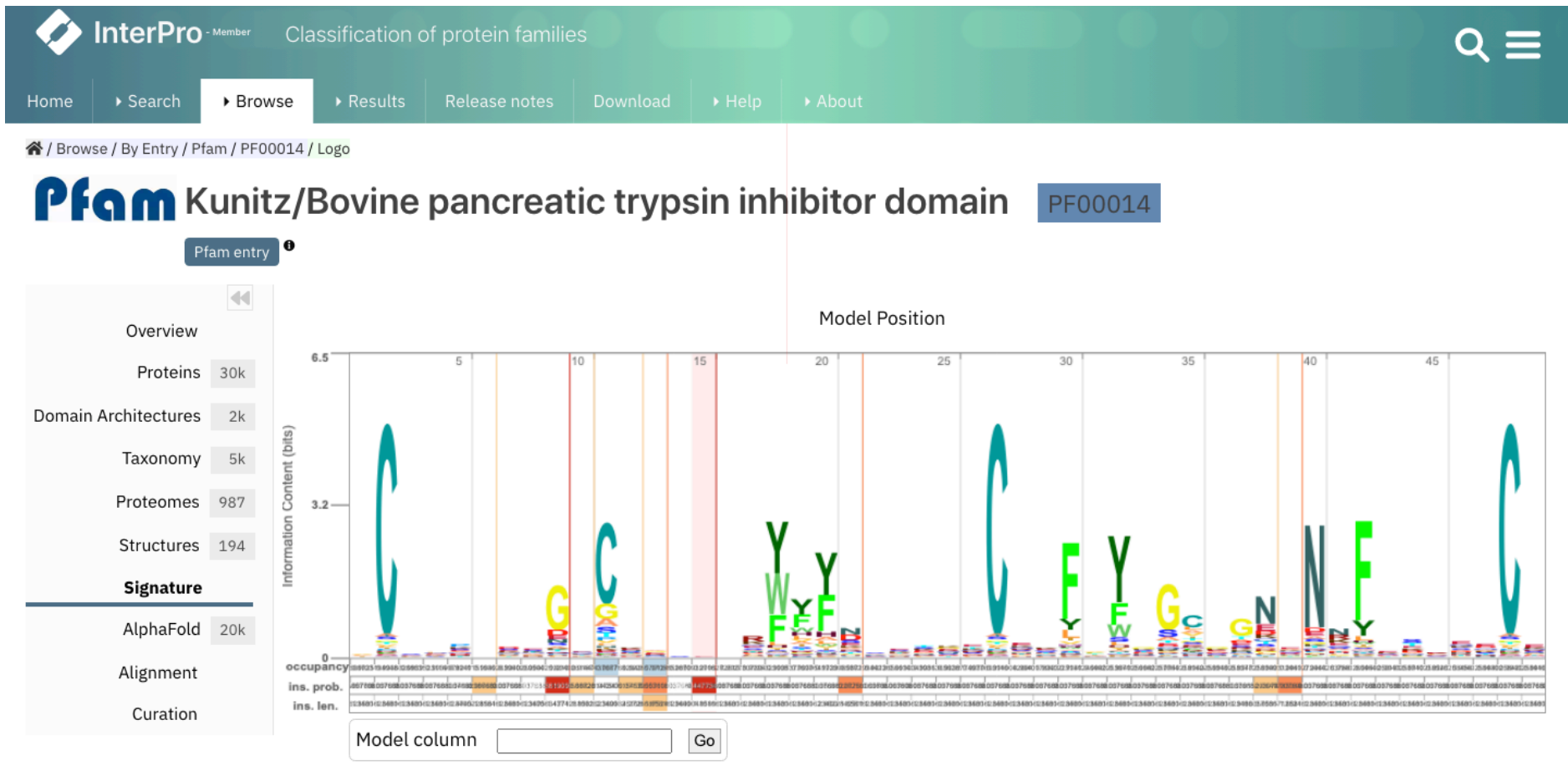
Author	Fenech M
Sequence Ontology	SO:0000417

HMM Information

HMM build commands	Build method: <code>hmmbuild -o /dev/null HMM SEED</code> Search method: <code>hmmsearch -Z 61295632 -E 1000 --cpu 4 HMM pfamseq</code>
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



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

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

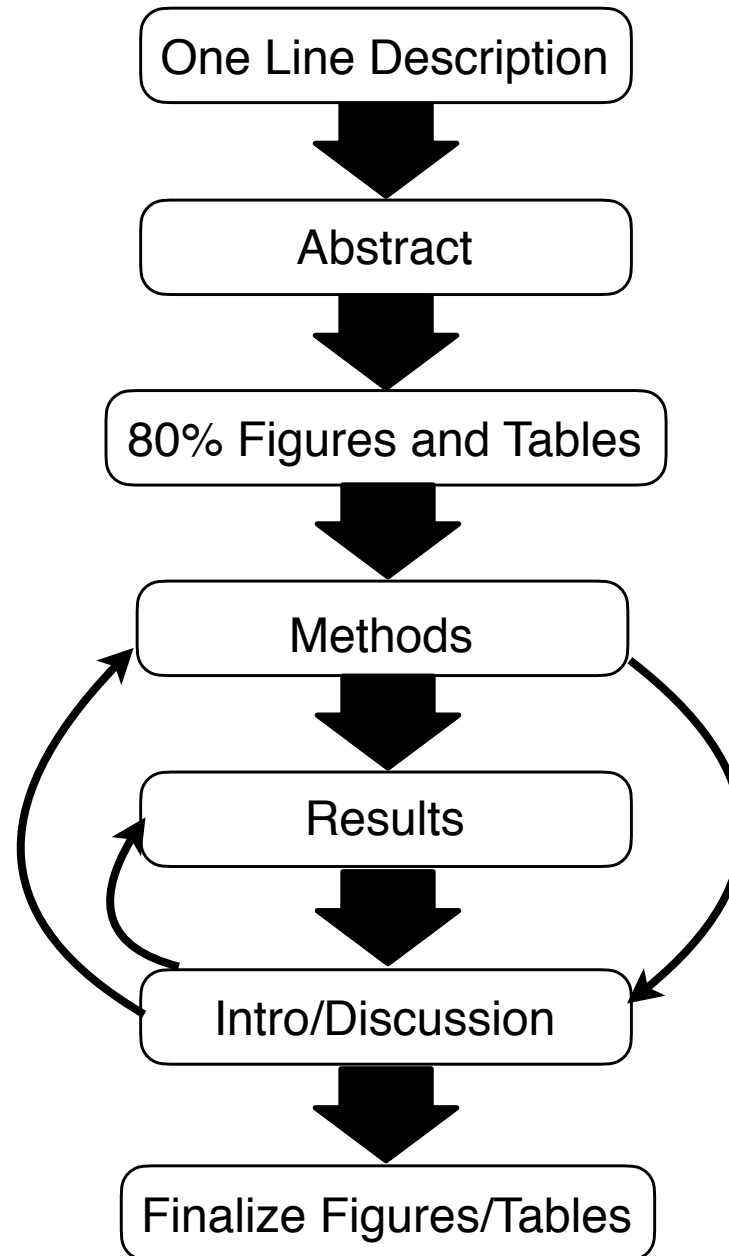
References

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 with subject: project-lb1b - Name Surname to:
emidio.capriotti@unibo.it

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**.