Анализ аминокислотной последовательности:

паттерны, домены, семейства

ИЛИ

что, где и как искать?

Что будем искать?







Основные понятия и термины

- Mecтo, <u>caйт</u>(site) -
- Momue (motif) –
- <u>Домен</u> (domain) –
- Семейство –
- Суперсемейство -

- Паттерн (pattern) –
- Позиционно специфическая матрица весов (PSSM) —
- Профиль-PSSM -
- Профиль-НММ -
- Подпись (signature) –
- «Отпечатки пальцев» (fingerprints) —
- Кластер -



эволюции, структуры и функции белков.

Домен – компактная, относительно независимо сворачивающаяся структура, относительно консервативная в процессе эволюции.

Белки могут состоять из одного или



Мотив?

- Мотив в аминокислотной последовательности набор консервативных остатков, важных для функции белка и расположенных на определенном (обычно коротком) расстоянии друг от друга в последовательности.
- Мотив структуры (структурный мотив) часто встречающийся в белках элемент пространственной структуры (α-спираль, β-шпилька, β-поворот).

В общем случае, структурные мотивы не обязательно соответствуют мотивам в аминокислотным последовательностях.

Один домен может содержать один или несколько мотивов в аминокислотной последовательности. Мотив может не входить в домены.

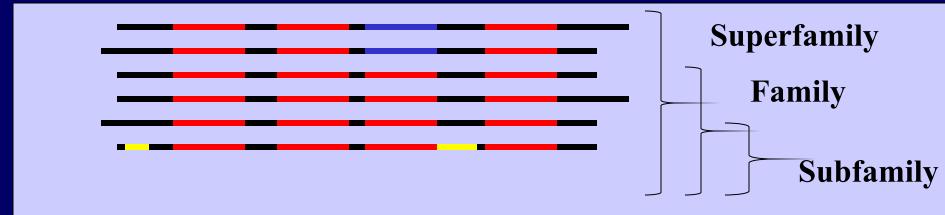
Не в любом выравнивании легко найти мотив.

Интуитивно понятно:

• Семейство - группа белков, имеющая общее происхождение, их аминокислотные последовательности выравниваются по всей длине со значимым весом и имеют сходную доменную структуру.

Мнения расходятся, когда речь идет о критериях:

• насколько должны быть похожи белки одного семейства (id>=30%, id>= 50%) ??? должны белки одного семейства выполнять одну и ту же функцию??



No comments

Основные понятия и термины

- Место, <u>сайт</u>(site) -
- Momue (motif) –
- <u>Домен</u> (domain) –
- Семейство –
- Суперсемейство -

- Паттерн (pattern) –
- Позиционно специфическая матрица весов (PSSM) —
- Профиль-PSSM -
- Профиль-НММ -
- Подпись (signature) –
- «Отпечатки пальцев» (fingerprints) -



InterPro

















The Protein Domain Database ProDom





Blocks WWW Server



Банки белковых семейств и доменов, производные от банков аминокислотных последовательностей

Коллекции мотивов

Коллекции доменов

PROSITE, 1989 BLOCKS

PRINTS

Pfam

SMART

ProDom, 1995

SUPERFAMILY

InterPro, 1999

(Integrated Resource of Protein Families)

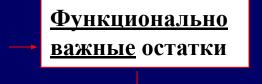
РКОSITE - <u>биологически</u> значимые сайты, паттерны и профили



Поиск

в SP

Выравнивание хорошо изученного семейства



4-5 консервативных остатков

Паттерн

Если находим только«правильные», то ОК

Если много лишнего, то увеличиваем паттерн

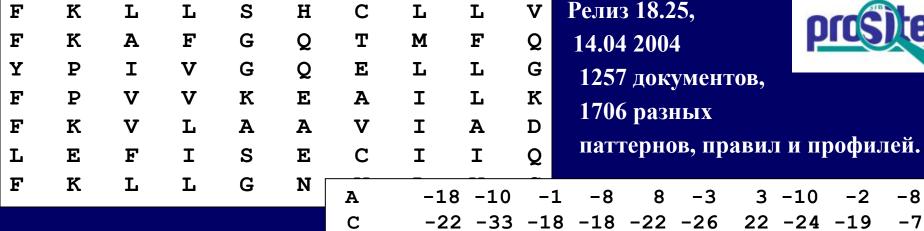
Паттерн – регулярное выражение UNIX'a:

$$[AC]-x-V-x(4)-\{ED\}$$

Ala или Cys- x-Val- x- x- x - x- (любой, но не Glu или Asp)

PROSITE - биологически значимые сайты, паттерны и профили

How we develop Prosite patterns!



Профиль или весовая матрица

-8 22 -24 -19 -7 -350 - 32 - 33-7 -17 -34 -31D 6 0 15 -25 -26 23 -27 -9 -9 -24 -23 E 14 -26 -29 -15 F 60 -30 12 4 12 - 29-30 -20 -28 -32 28 -14 -23 -33 -27 G -13 -12 -25 -25 -16 14 -22 -22 -23 -10 Η I 3 - 2721 25 -29 -23 -8 33 19 -23 25 -25 -27 K -6 4 -15 -27 -26 L 14 -28 19 27 -27 -20 -9 33 26 -21 M 3 -15 10 14 -17 -10 -9 25 -22 -6 -24 -278 -15 -24 -24 N 1 -30 24 -26 -28 -14 -10 -22 -24 -26 -18 P 5 -25 -26 24 -16 -17 -23 Q -32 -9 -18 9 -22 -22 -10 0 -18 -23 -22 R S -22 -8 -16 -21 11 2 -1 -24 -19Т -10 -10 -6 -7 -5 -8 2 -10 V -25 22 25 -19 -26 6 19 16 -16 9 -25 -18 -19 -25 -27 -34 -20 -17 -28 W 34 -18 -1 1 -23 -12 -19 0 Y 0 - 18

Pfam



- http://www.sanger.ac.uk/Software/Pfam/index.shtml
- Большая коллекция множественных выравниваний, доменов, семейств и профилей-НММ для них.
- Состоит из 2-х частей:
 - PfamA курируемая часть, покрывает 73% SWISS-Prot+TrEMBL
 - PfamB большое число маленьких семейств из автоматически сгенерированной базы доменов ProDom, не вошедших в PfamA.
- Удобна для анализа доменной структуры белков.



Pfam



- 1. Множественное выравнивание (ClustalX) некоторого семейства или кластера.
- 2. Экспертиза и корректировка выравнивания- затравки.
- 3. Построение профиля-НММ для затравки.
- 4. Поиск в базе данных а.к.последовательностей новых членов данной группы.

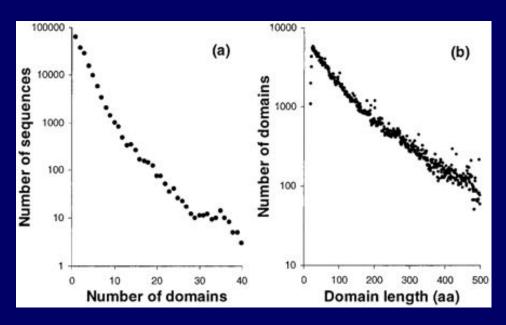
ProDom

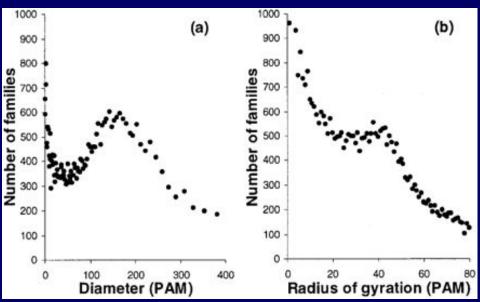


- http://www.toulouse.inra.fr/prodom.html
- Рассматриваются все последовательности в SWISS-Prot+TrEMBL.
- Автоматическое выделение доменов (программа DOMAINER: сначала локальное попарное выравнивание (blastp) всех против всех, затем кластеризация)
- Коллекция доменов >150 000 семейств.
- Некоторые семейства выделены на основе выравниваний из PfamA.
- Гомогенность семейства оценивается с помощью диаметра (тах расстояния между 2 доменами в семействе) и радиуса (ср.кв. расстояние между доменами и консенсусом семейства). Оба параметра измеряются в РАМ

Статистика ProDom

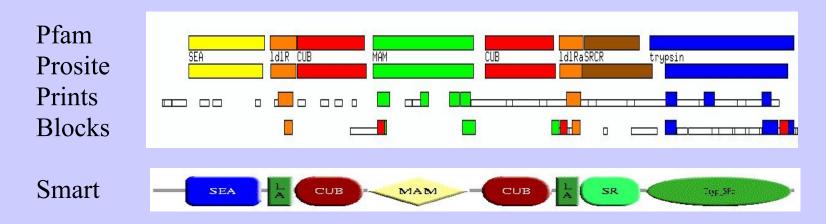






Всего — 157 167 семейств.
43 965 из них содержат
более 2 последовательностей.
Среднее число доменов в
последовательности — 2.8
Средняя длина — ~ 130
а.к. остатков

Comparison of protein family databases: an example

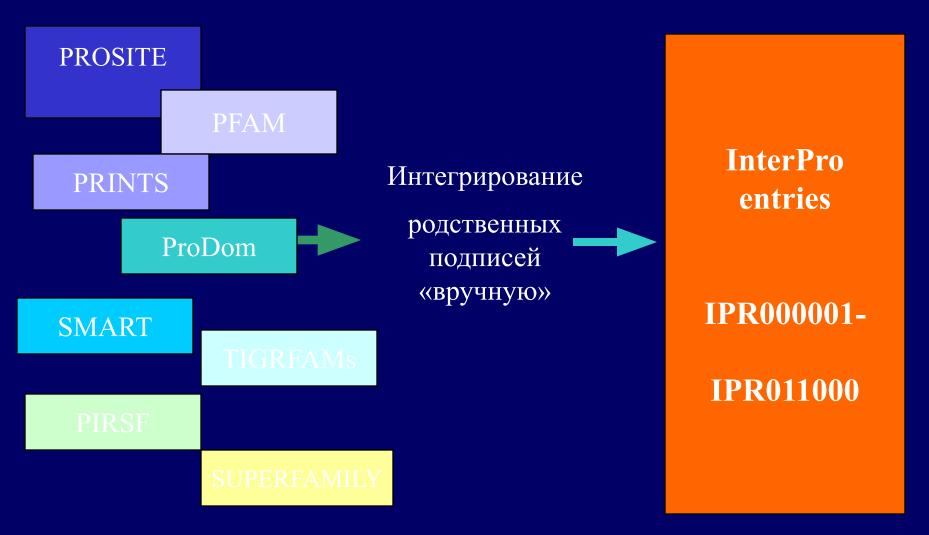


(ProDom, PIRaln, ProClass, Systers, Picasso etc. not shown)

Example: ENTK_HUMAN (Enteropeptidase precursor)

Создание интегрированной базы данных InterPro





InterPro- an <u>integrated</u> resource of <u>pro</u>tein families, domains and functional sites.

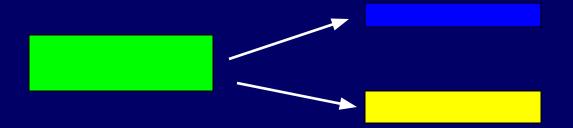


Entry types in InterPro

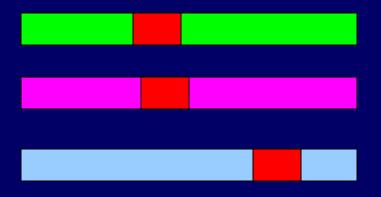
- Family group of evolutionarily related proteins, that share one or more domains/repeats in common.
- **Domain** -independent structural unit which can be found alone or in conjunction with other domains or repeats.
- Repeat -region occurring more than once that is not expected to fold into a globular domain on its own.
- **PTM** (post-translational modification) -The sequence motif is defined by the molecular recognition of this region in a cell.
- Active site -catalytic pockets of enzymes where the catalytic residues are known.
- **Binding site** —binds compounds but is not necessarily involved in catalysis.

Взаимосвязи подписей в InterPro

• Parent/child — уровень семейства

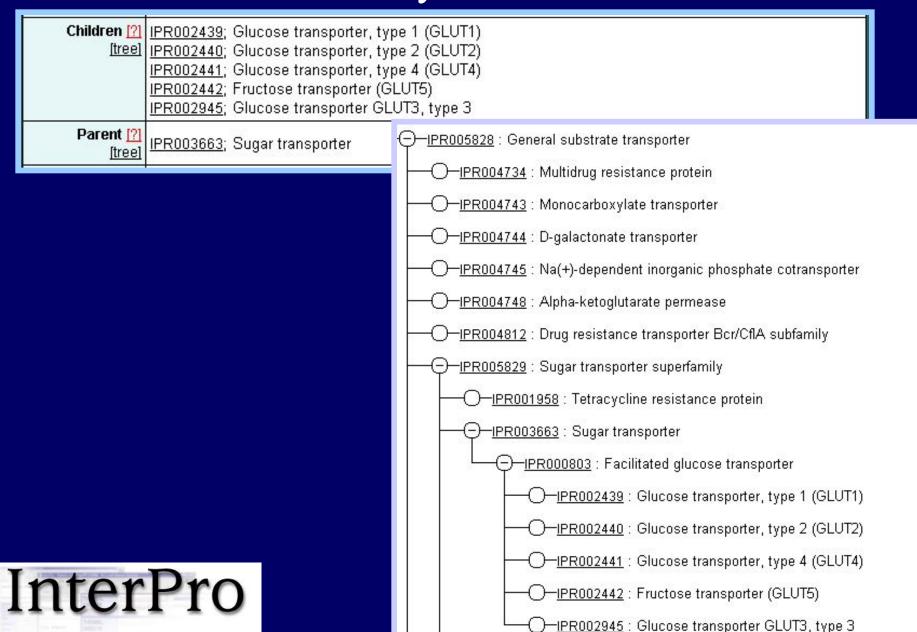


• Contains/found in ___ состав домена



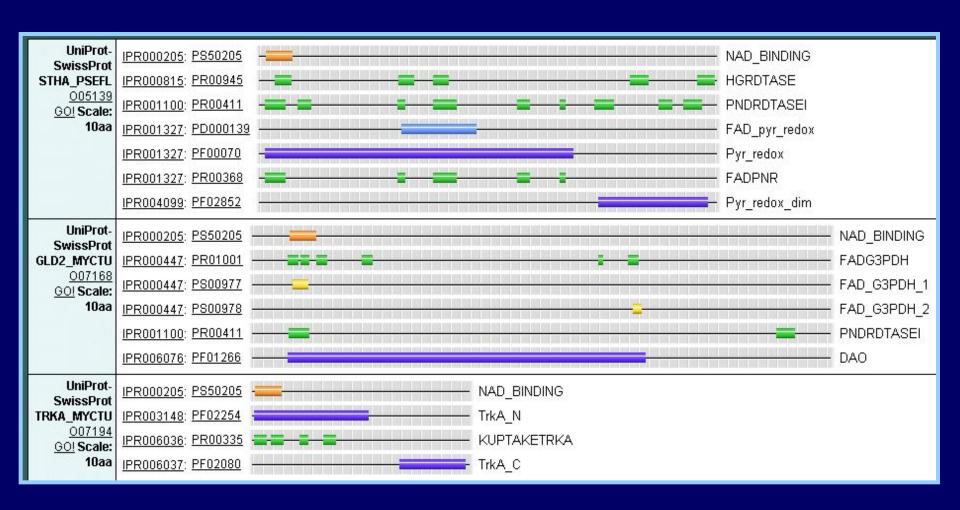
InterPro

Parent/child- family level





Contains/found in



PROTOMAP

ProtoMag

- http://www.protomap.cs.huji.ac.il
- Automatic classification of all SWISS-PROT proteins into groups of related proteins (also including TrEMBL now)
- Based on pairwise similarities
- Has hierarchical organisation for sub- and super-family distinctions
- 13 354 clusters, $5869 \ge 2$ proteins, $1403 \ge 10$
- Keeps SP annotation eg description, keywords
- Can search with a sequence -classify it into existing clusters

ATPO_P TG

ersion 3.0

You are in 'Simplified mode

