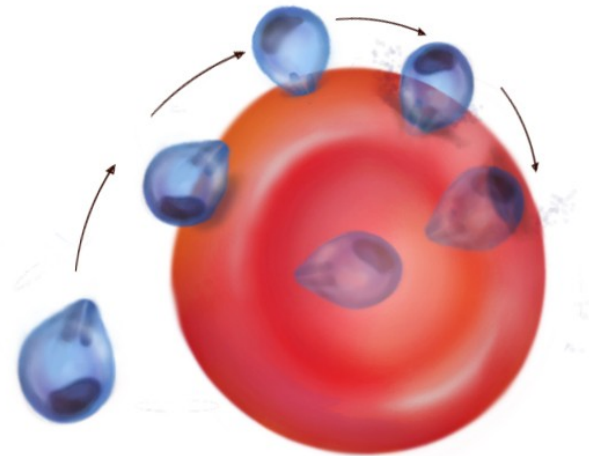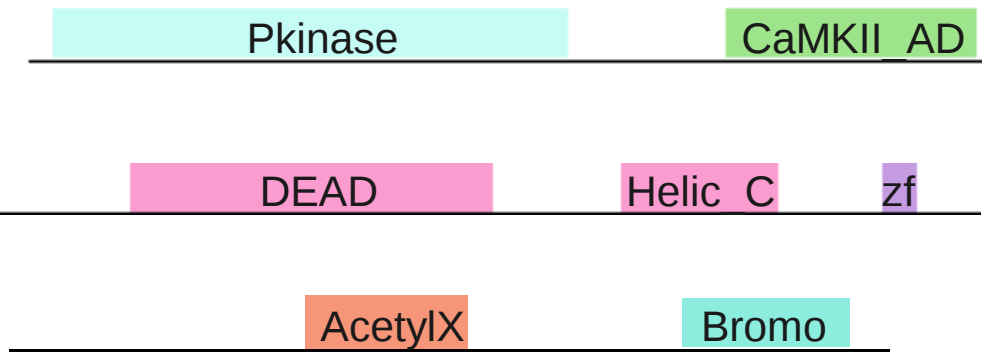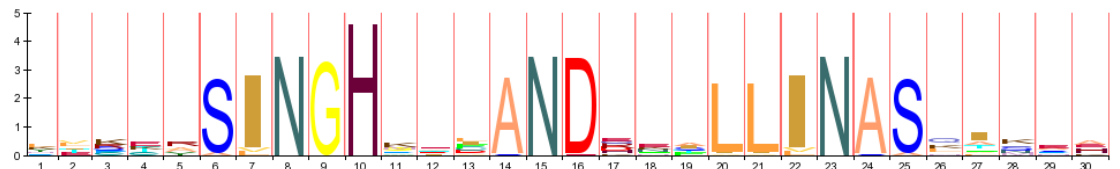# Forget the *E*-value: family-based *q*-values for protein domain prediction, and empirical error detection
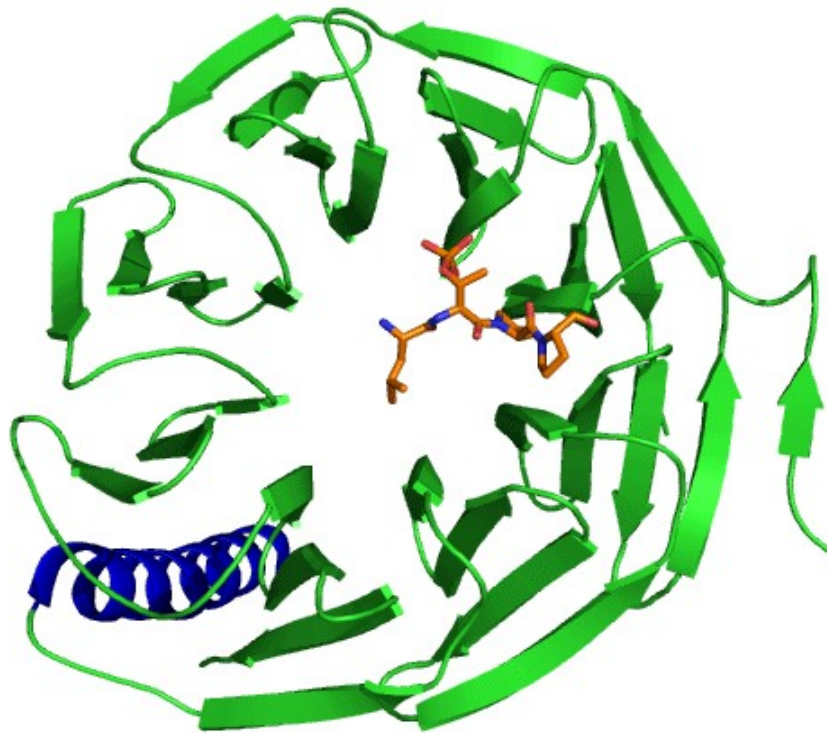
Alejandro Ochoa

Molecular Biology, Princeton University
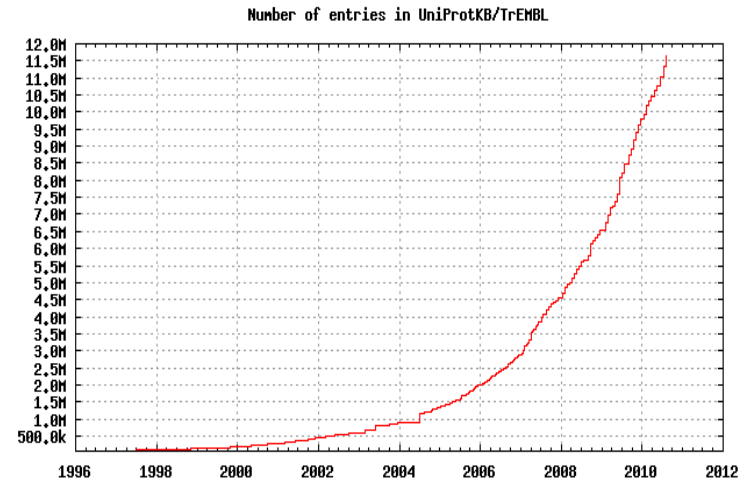
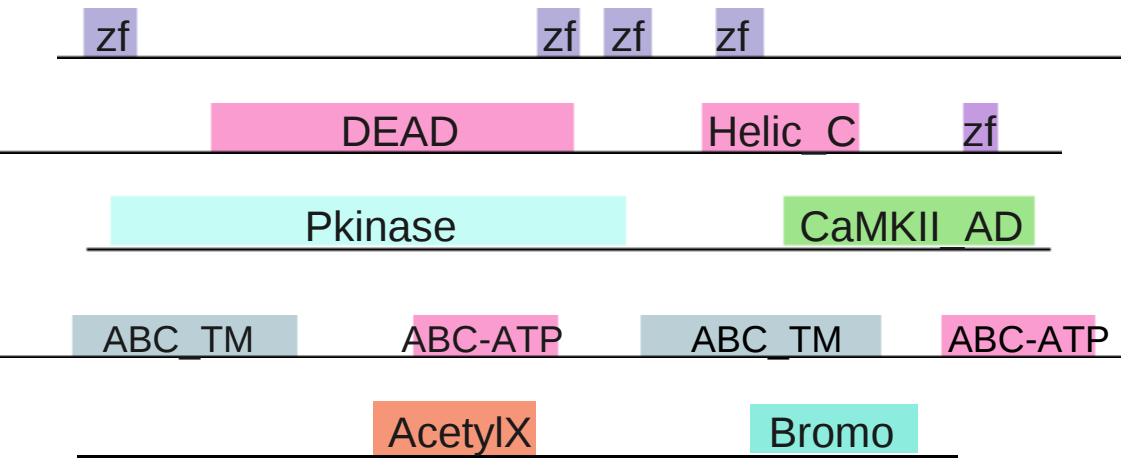NCBI, NIH, 2013-02-25

# Protein domains



Structure

Evolution

Function

Sequence-based domain prediction:

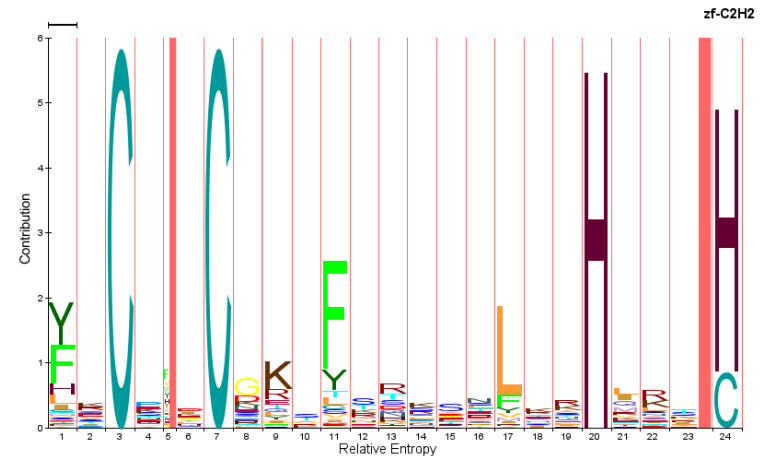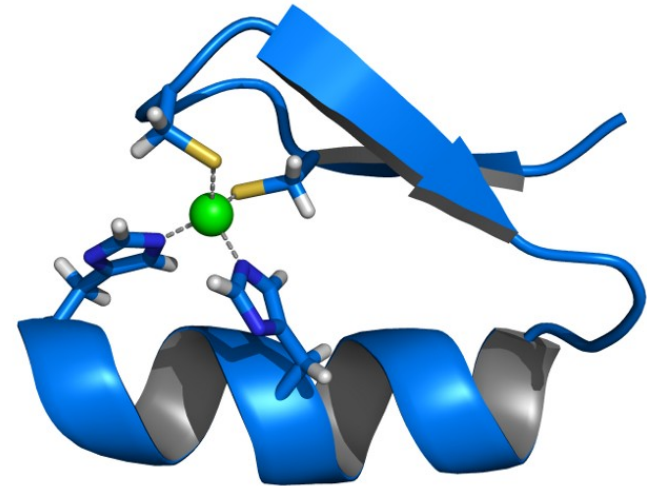# Why predict domains?



For new sequences, before experiments start…

Domains may imply functions

Experimental alternatives are unfeasible as protein databases grow exponentially

# Representing Domains

# Databases of Domain Families

This work uses Pfam and HMMER, but theory and results are general

Two members of Pfam "Clan" HTH

# *Plasmodium falciparum*

Malaria

Information challenges

- Diverged eukaryote
- 80% AT-bias
- Low-complexity regions

Annotation

- 5.5K proteins
- 45% unknown function
  - 20% unknown in yeast
- 88% of annotations are bioinformatical

# Sequence-based annotation of *Plasmodium falciparum* proteome

# Poor domain coverage of *Plasmodium falciparum*



- Known domain families, very diverged
- Novel domain families
- Disordered sequences

# Outline of results

- <span style="color:red">Domain prediction using context</span>
  - Application to malaria parasite
- Optimal *FDR* control for domains
  - Family-based *q*-values, adapted to domain problem
  - *Local FDR* optimizes problem, but *q*-values are more robust
  - Large problems with coiled coils, transmembrane domains
  - Implications for sequence models

# Domain Prediction Using Context: dPUC



DEAD    Helic_C    zf

RNA Pol II domains

## Background

– Domains co-occur in limited combinations

– Domains are scored independently of each other

## Idea

– Score domains in combination

– Context + Sequence evidence

# The dPUC method

**A**



**B**



**C**

**Standard Pfam:**



**dPUC Pfam:**



http://compbio.cs.princeton.edu/dpuc/
Ochoa, *et al*. BMC Bioinformatics 12, 90 (2011).

# Improved signal to noise

**Real protein**

**Real protein with shuffled sequences**



*P. falciparum*

- ■ **Pfam Standard**
- ✕ **Pfam Extended**
- ○ **dPUC**
- ▢ **CODD**

http://compbio.cs.princeton.edu/dpuc/
Ochoa, *et al*. BMC Bioinformatics 12, 90 (2011).

# Improved ortholog coherence on *Plasmodium* species

# New predictions

Phosphatase -> RNA lariat debranching enzyme

*P. falciparum*

Std Pfam

dPUC

Metallophos          DBR1

*S. cerevisiae*

Std Pfam & dPUC

# New predictions

MIF4G domain-containing protein ->
  Poly-A binding protein-interacting protein 1

## *P. falciparum*



## *H. sapiens*

Interested malaria curator:
Hagai Ginsburg,
Hebrew U of Jerusalem

# New predictions

## RNA helicase -> mRNA sequestration

# dPUC increases coverage

# Domain context

Pkinase    CaMKII_AD

DEAD    Helic_C    zf

AcetylX    Bromo

Complements sequence evidence

Improves domain predictions

Works best on diverged organisms

# The dPUC method: problems



Domain scores are normalized by curated thresholds

Why do *E*-values perform worse?

# Outline of results

- Domain prediction using context
  - Application to malaria parasite
- <span style="color:red">Optimal *FDR* control for domains</span>
  - Family-based *q*-values, adapted to domain problem
  - *Local FDR* optimizes problem, but *q*-values are more robust
  - Large problems with coiled coils, transmembrane domains
  - Implications for sequence models

# The *E*-value: the statistic of sequence analysis

**Query Type**  **Database Type**  **Software**

GVYFEKSRSSWTA    GVRQRKNSNKWVS    **FASTA (*Z* 1985)**
**BLAST (*E* 1990)**



GVRQRKNSNKWVS    **PSI-BLAST (*E* 1997)**

GVYFEKSRSSWTA



**SAM (*Z* 1994, ~*E* 2005)**
**HMMer (bit 1995, ~*E* 1998, *E* 2008)**
**Pfam DB (curated 1997)**
**IMPALA, RPS-BLAST (*E* 1999)**

    

**PRC (bit 2002, ~*E* 2004)**
**HHsearch (~*E*,*P* 2004)**
**CORAL (*E* 2009)**

*FDR*: **BH 1995, *q* 2002**

# FDR = E/n

*FDR* = average posterior error probability

# Computing *q*-values

*P*-values of random data are uniformly distributed

*P*-values of real data has two components

(when *p*-values are accurate)

1

$\pi_0$

0                    p                    1      0                    p                    1

Step 1: estimate $\pi_0$ (proportion of data that is false)

# Computing *q*-values

$$A_{\mathrm{null}}(t) = \pi_0 t,$$

$$A_{\mathrm{all}}(t) = \frac{n(t)}{N},$$

$$\mathrm{FDR}(t) = \frac{A_{\mathrm{null}}}{A_{\mathrm{all}}} = \frac{\pi_0 t N}{n(t)},$$

$\pi_0$

$\pi_0 t$

0   $t$   **p**   1

Step 2: Directly
  estimate *FDR(t)* for
  all thresholds *t*
   – *N* = # tests
   – *n(t)* = # sig tests

Step 3: Ensure
  monotonicity

$q(p) = \min_{t;p \le t} FDR(t)$

# *FDR* and *local FDR*

$$A_{\text{null}}(t) = \pi_0 t, \qquad h_{\text{null}}(t) = \pi_0,$$

$$A_{\text{all}}(t) = \frac{n(t)}{N}, \qquad h_{\text{all}}(t) = \frac{n'(t)}{N},$$

$$\text{FDR}(t) = \frac{A_{\text{null}}}{A_{\text{all}}} = \frac{\pi_0 t N}{n(t)},$$

$$\text{local FDR}(t) = \frac{h_{\text{null}}}{h_{\text{all}}} = \frac{\pi_0 N}{n'(t)}.$$

$\pi_0$

$\pi_0 t$

0   $t$   p   1

*Local FDR* =
  Posterior Error
  Probability (*PEP*)

*FDR* = average *PEP*
  of significant
  predictions

# *Local FDRs* optimize domain prediction

Find domain family thresholds $t_i$ (for each family $i$) to maximize predictions

$$M = \sum n_i\left(t_i\right),$$

while constraining the combined *FDR* of all families to *Q*

$$Q \geq \frac{\sum \pi_{0,i} \cdot t_i \cdot N_i}{\sum n_i\left(t_i\right)} = \frac{\sum FDR_i\left(t_i\right) \cdot n_i\left(t_i\right)}{\sum n_i\left(t_i\right)}.$$

Equal family *local FDR*s solve this optimization!

# *Q*-values for domains

Introduced domain-specific features

- – Incomplete list of *p*-values (HMMER3, BLAST and descendants)
- – Correction for domain overlaps (heuristic)

Two members of Pfam "Clan" HTH

$$FDR = E/n$$

# Pfam curated thresholds

# Empirical null models
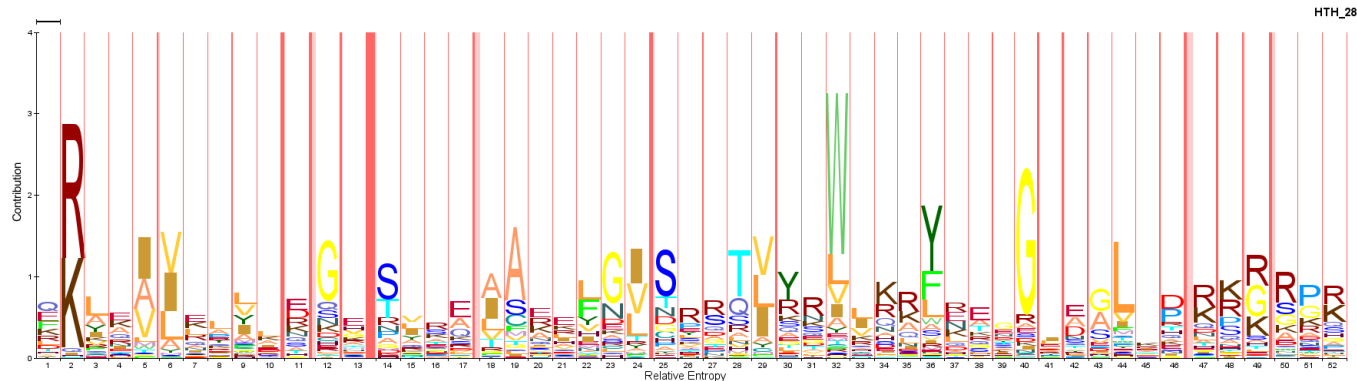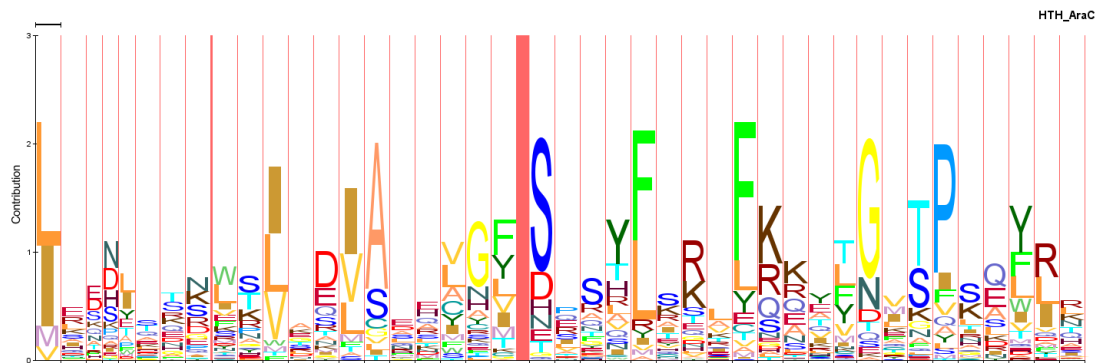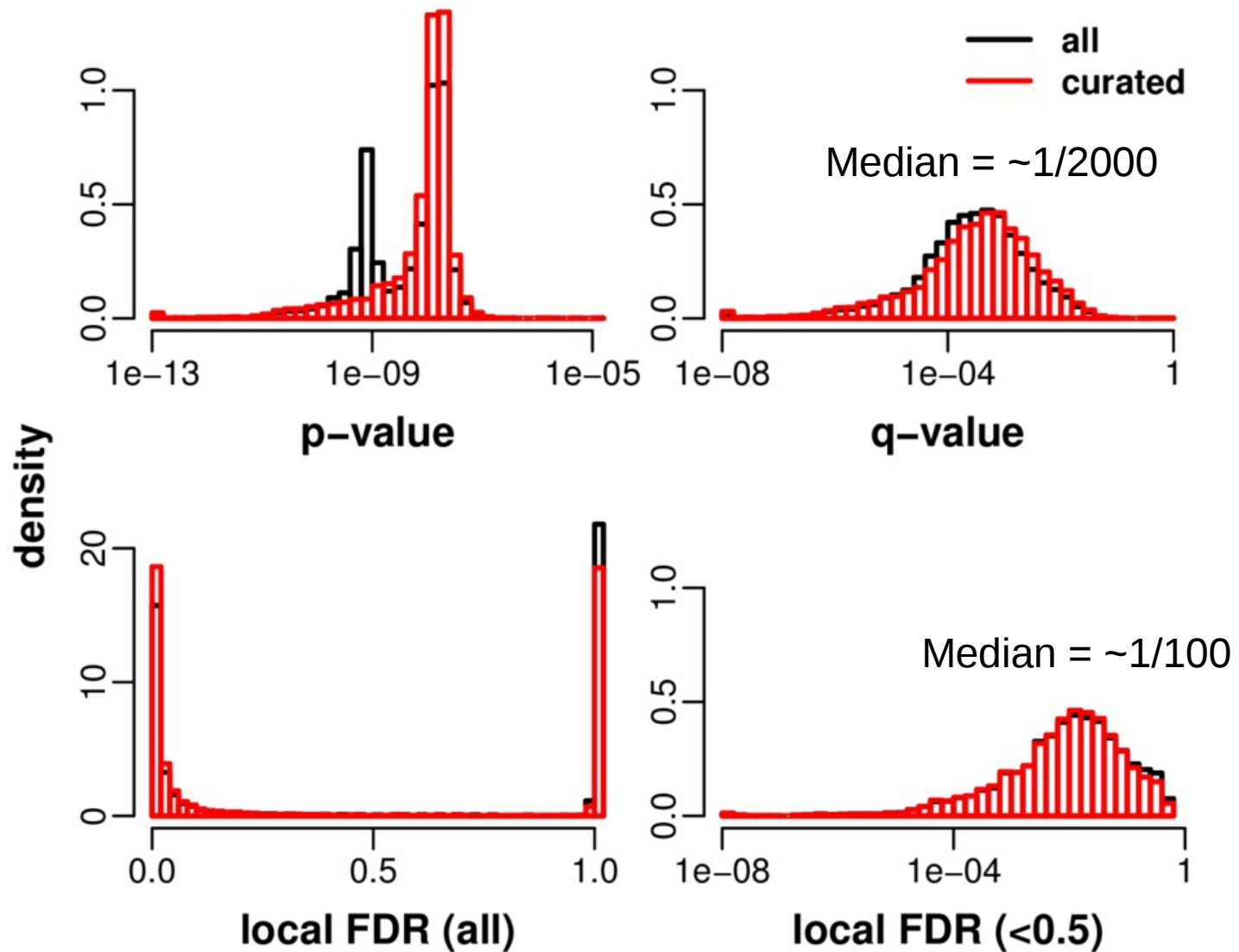
$$FDR = \#FP / ( \#FP + \#TP )$$

## 2nd order Markov Random Sequences

- **FP if domain came from random sequence**

    Improved from Ochoa, *et al*. BMC Bioinformatics 12, 90 (2011).

## Ortholog Set Coherence

- **FP if orthologs don't predict any homologous domains**

    Improved from Ochoa, *et al*. BMC Bioinformatics 12, 90 (2011).

## Clan Overlap

- **FP if domain overlaps stronger non-homologous domain**

    Inspired/adapted from S. Eddy (p.c., 2012).
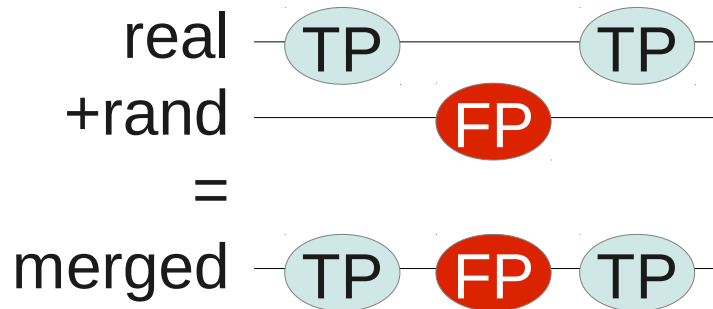
## Context Coherence

- **FP if domain doesn't co-occur with any stronger domains**

    Inspired/adapted from Terrapon, *et al*. BMC Bioinformatics 13, 67 (2012).

# Empirical null models

$$FDR = \#FP \, / \, ( \, \#FP + \#TP \, )$$

**2$^{nd}$ order Markov Random**

real   TP  —  TP

+rand  —  FP  —

=

merged   TP — FP — TP

**Ortholog Set Coherence**

TP —— TP

orthologs
<50% ID

FP TP TP

same
clan

**Clan Overlap**

1 seq,
keep
overlaps

TP   TP   TP

Ø   FP   Ø

Ø

same
(diff)
clan

**Context Coherence**

1 seq   TP — FP — TP

(un)observed
clan pair

# *q > local FDR*

# Improving more than just domains

# The *q*-values underestimate empirical *FDR*s

# Per family analysis of noise

# Classes of noisy domains



**Coiled coils: 403 (3%)**

**Transmembranes: 1466 (12%)**

**Low complexity: 1987 (16%)**

**Cysteine-rich: 587 (5%)**

**Normal: 7830 (64%)**

Similar but non-homologous families (different clans)

Rackham *et al*. JMB 403, 480–493 (2010).
Wong, *et al.* PLoS Comput Biol 6, e1000867 (2010).

# Classes enriched in noisy domains

# *Local FDR > q*
# in families with correct stats



**Domain increase (%)** (y-axis)

**Empirical FDR** (x-axis)

**2° Markov Random (MR2)**

- local FDR
- q-value
- Standard Pfam
- E-value

**Ortho Set Coher (OSC)**

**Clan Overlap (CO)**

**Context Coherence (CC)**

# Conclusions

The *q*-value and *local FDR* are better for domains

- And likely better for regular sequence database searches, iterated searches, orthology prediction

- *E*-values do not control posterior error probability

Presented novel empirical null models

- Needed to verify theory is correct

- Uses common-sense biological information and real, full protein sequences

- Structural benchmarks (i.e. SCOP) are limited to well-studied, single domains from model organisms, and exclude coiled coils and transmembrane domains

# Noisy domains



**Coiled coils**                    **Transmembranes**

These remain large problems in all sequence analysis.

Solutions?

- Cannot prevent or ignore these queries

- Masking removes too much information

- Benchmarks not powerful enough to give better thresholds

- Can we properly handle these common, correlated patterns that do not imply homology?

Rackham *et al*. JMB 403, 480–493 (2010).
Wong, *et al.* PLoS Comput Biol 6, e1000867 (2010).

**PRINCETON UNIVERSITY**

- **Mona Singh, Computer Science**
  - Jesse Farnham
  - Dario Ghersi
  - Peng Jiang
  - Shilpa Nadimpalli
  - Anton Persikov
  - Yuri Pritykin
  - Pawel Przytycki
  - Josh Wetzel

- **Thesis Committee**
  - John Storey
  - Leonid Krugliak
  - Saeed Tavazoie

- **Manuel Llinás, Molecular Biology**
  - Lindsey Altenhofen
  - Katie Baska
  - Simon Cobbold
  - Valerie Crowley
  - Björn Kafsack
  - Ian Lewis
  - Jessica O'Hara
  - Heather Painter
  - Joana Santos
  - Ariel Schieler
  - April Williams

- **NSF GRFP**

http://compbio.cs.princeton.edu/dpuc/