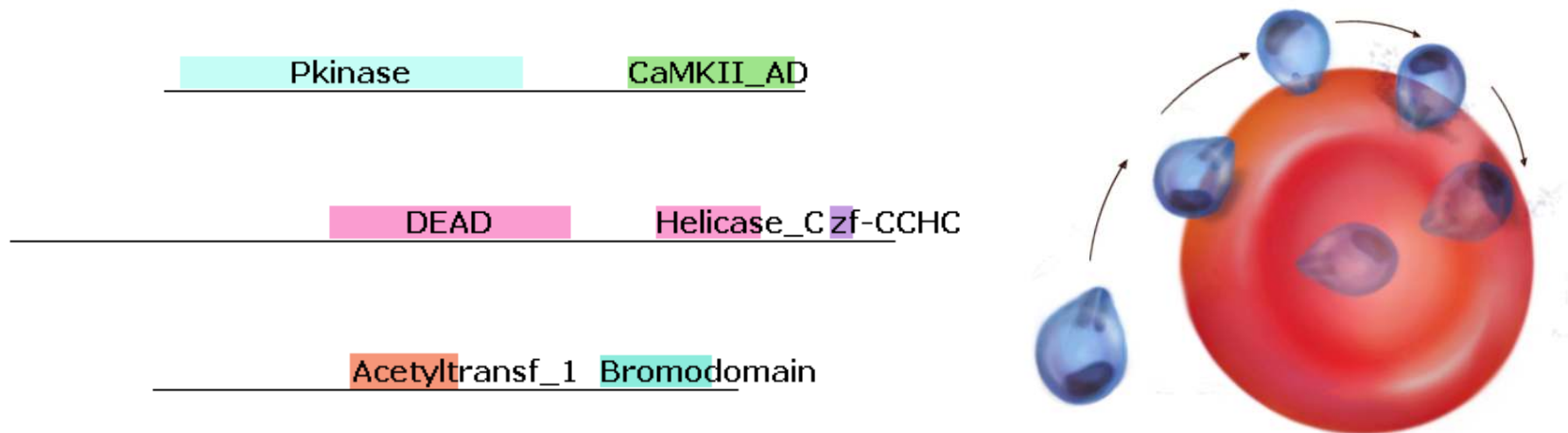
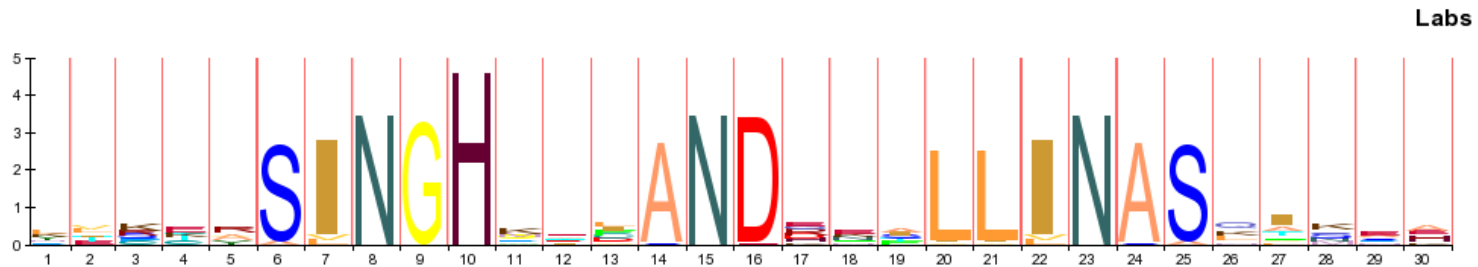


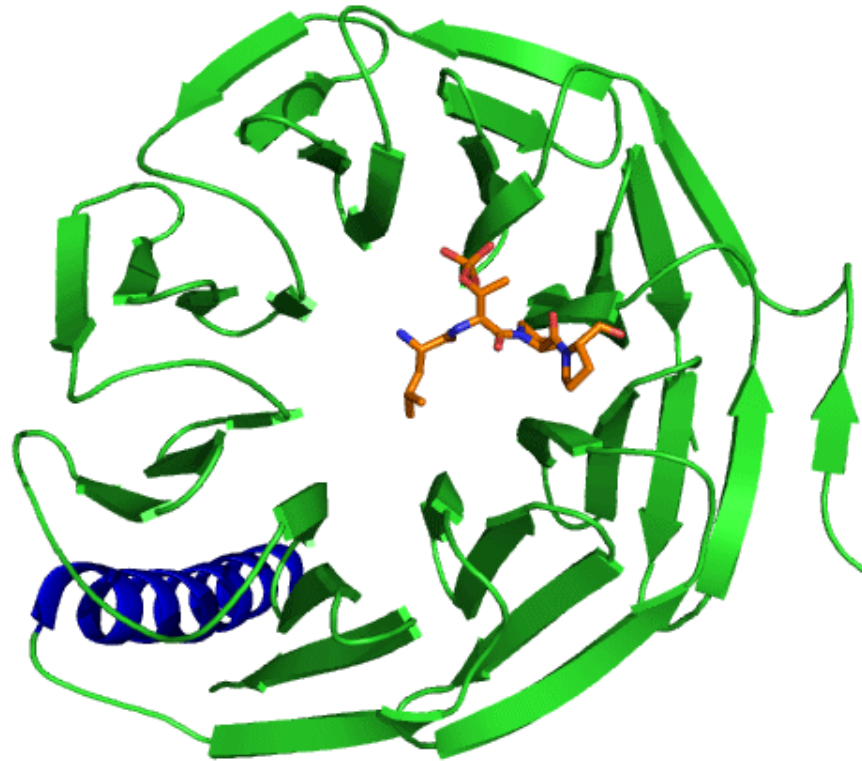
# Improving Domain Prediction in *Plasmodium falciparum*



Alejandro Ochoa  
2010-10-09



# Protein domains



Domain predictions:

F-box

WD4 WD40 WD40

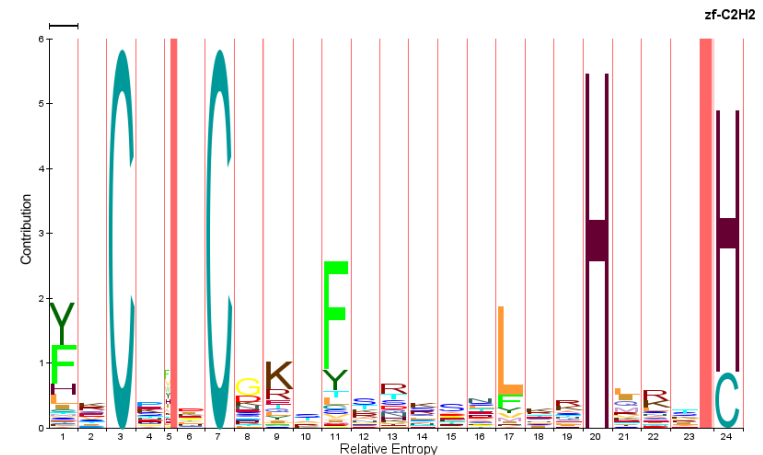
WD40 WD40 WD40

# Pfam: a database of protein domain families

SNAI\_DROME/362-385  
 SNAI\_XENLA/232-255  
 SNAI\_MOUSE/236-259  
 ESCA\_DROME/426-449  
 SUHW\_DROAN/221-243  
 TERM\_DROME/323-346  
 Z020\_XENLA/174-196  
 EVI1\_HUMAN/217-239  
 Z02\_XENLA/34-59  
 EVI1\_HUMAN/21-44  
 ZNF10\_HUMAN/517-539  
 ZNF91\_HUMAN/238-260  
 ZFP58\_MOUSE/120-142  
 TRAI\_CAEEL/306-331  
 ZNF76\_HUMAN/345-368  
 ZN12\_MICSA/106-129  
 LOLAI\_DROME/794-817  
 ZNF17\_HUMAN/435-457  
 ZG32\_XENLA/34-56  
 TF3A\_BUFAM/104-128  
 ZG46\_XENLA/146-168  
 MZFI\_HUMAN/412-434  
 ZN239\_MOUSE/6-28  
 ZSC22\_HUMAN/352-374  
 EGR1\_HUMAN/396-418  
 SUHW\_DROAN/349-373  
 CF2\_DROME/485-508  
 CF2\_DROME/401-423  
 KRUP\_DROME/306-328  
 TTY1\_HUMAN/383-407  
 ZG52\_XENLA/61-83  
 TTKB\_DROME/538-561  
 ZNF76\_HUMAN/285-309  
 SDC1\_CAEEL/145-168  
 SRYC\_DROME/358-380  
 SDC1\_CAEEL/270-292  
 TRAI\_CAEEL/276-300  
 ESCA\_DROME/370-392

YACQ...VCH...KSF SRM...SLLNKHSS...NC  
 YQCK...SCS...RTFSRM...SLLHKHEET...GC  
 YQCQ...ACA...RTFSRM...SLLHKHSES...GC  
 YSCT...SCS...KTFSRM...SLLTKHSEG...GC  
 HVCG...KCY...KTFRRL...MSLKKHLEF...C  
 LHCR...RCR...TQFSRR...SKLHIHQKL...RC  
 FMCA...DCG...RCFSVS...SSLKYHQRI...C  
 IKCK...DCG...QMFSTT...SSLNKHRRF...C  
 YSCA...DCG...KHFSEK...MYLQFHQKMPSEC  
 YRCE...DCD...QLFESK...AELADHQKF...PC  
 YKCN...QCG...IIFSQM...SPFFIVHQIA...H  
 YKCE...ECG...KAFKQL...STLTTHKII...C  
 IKCE...ECG...KAFSTR...STYYRHQKN...H  
 YKCEF.ADCE...KAFSNA...SDRAKHQNR...TH  
 YTCS...TCG...KTYRQT...STLAMHKRS...AH  
 YRCS...QCG...KAFRRT...SDLSSHRT...QC  
 YECR...HCG...KKYRWK...STLRRHENV...EC  
 YECN...KCG...KFFRYC...FTLNRHQRV...H  
 FVCV...HCG...KGFSDM...YKLSLHLRI...H  
 YVCYF.ADCG...QQFRKH...NQLKIHQYI...H  
 YVCT...ECG...TSFRVR...PQLRIHLRT...H  
 FVCG...DCG...QGFVRS...ARLEEHRV...H  
 YKCD...KCG...KGFTRS...SSLLVHHSV...H  
 YKCG...ECG...KTFSRS...THLTQHQRV...H  
 FACD...ICG...RKFARS...DERKRHTKI...H  
 YACK...ICG...KDFTRS...YHLKRHQKYS...SC  
 YTCP...YCD...KRFTQR...SALTVHTTK...LH  
 YTCS...YCG...KSFTQS...NTLKQHTRI...H  
 YTCE...ICD...GKFSDS...NQLKSHMLV...H  
 YVCPF.DGCN...KKFAQS...TNLKSHILT...H  
 YTCT...QCN...KQFSHS...AQLRAHIST...H  
 YPCP...FCF...KEFTRK...DMMTAHVKI...IH  
 YTCPE.PHCG...RGFTSA...TMYKNHVRI...H  
 YMCQ...VCL...TLFGHT...YNLFMHVRT...SC  
 YOCD...ICG...QKFVOK...INLTHHARI...H  
 YFCH...ICG...TVFIEQ...DMLFKHWRL...H  
 NKCEY.PGCG...KEYSRL...ENLKTHTRT...H  
 CKCN...LCG...KAFSRP...WLLQGHIRT...H

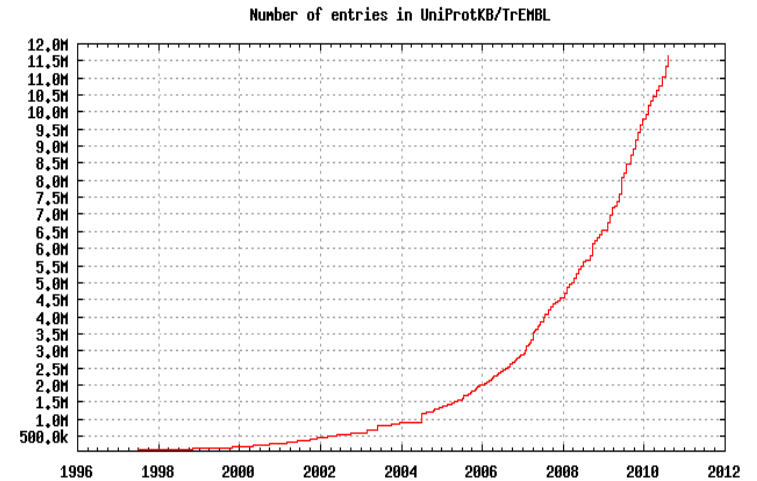
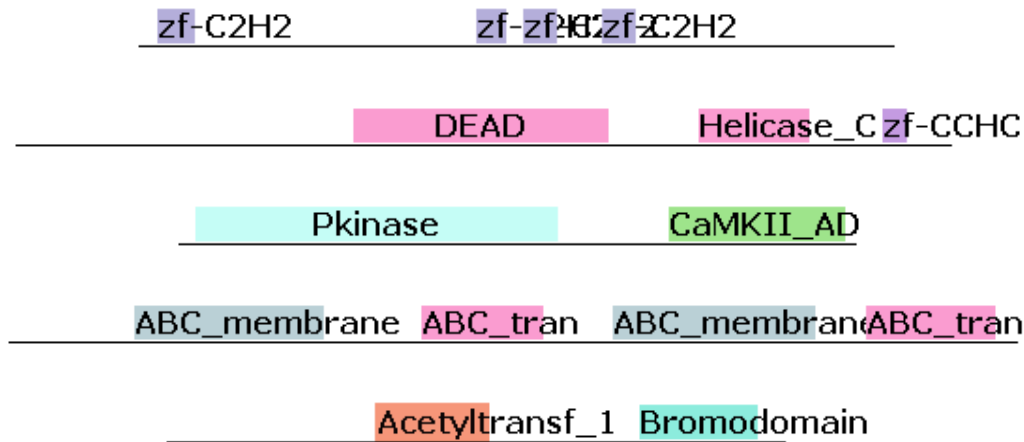
- 11,912 curated families!
- Profile Hidden Markov Models (HMMs): probabilistic models of sequence families



■ ■ ■

■ ■ ■

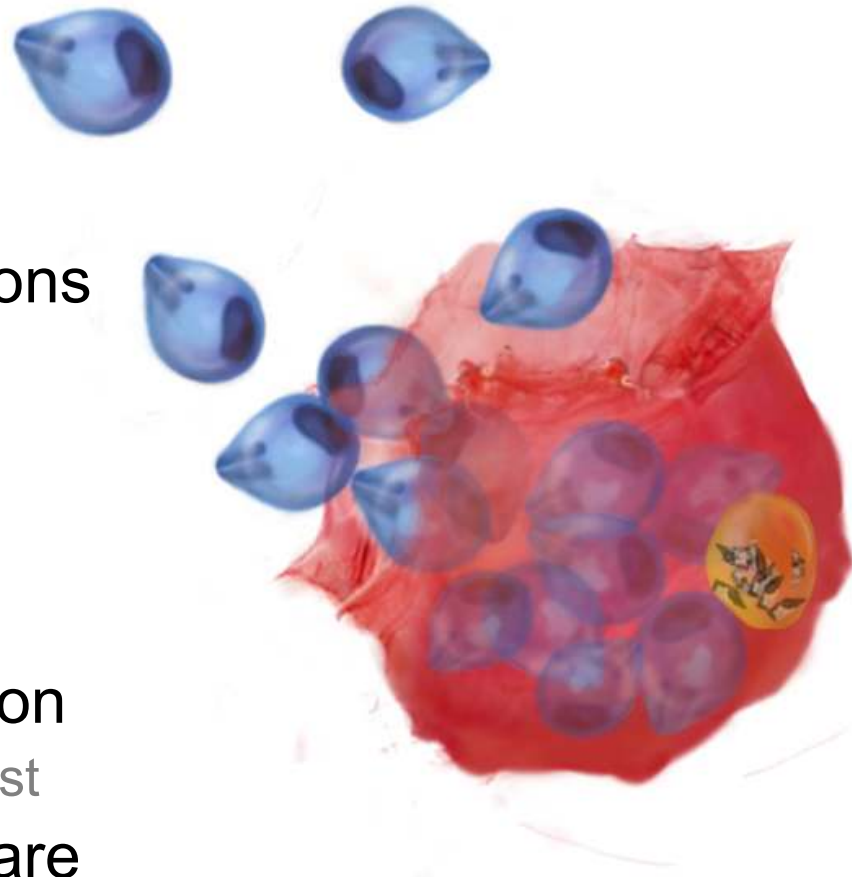
# Why predict domains?



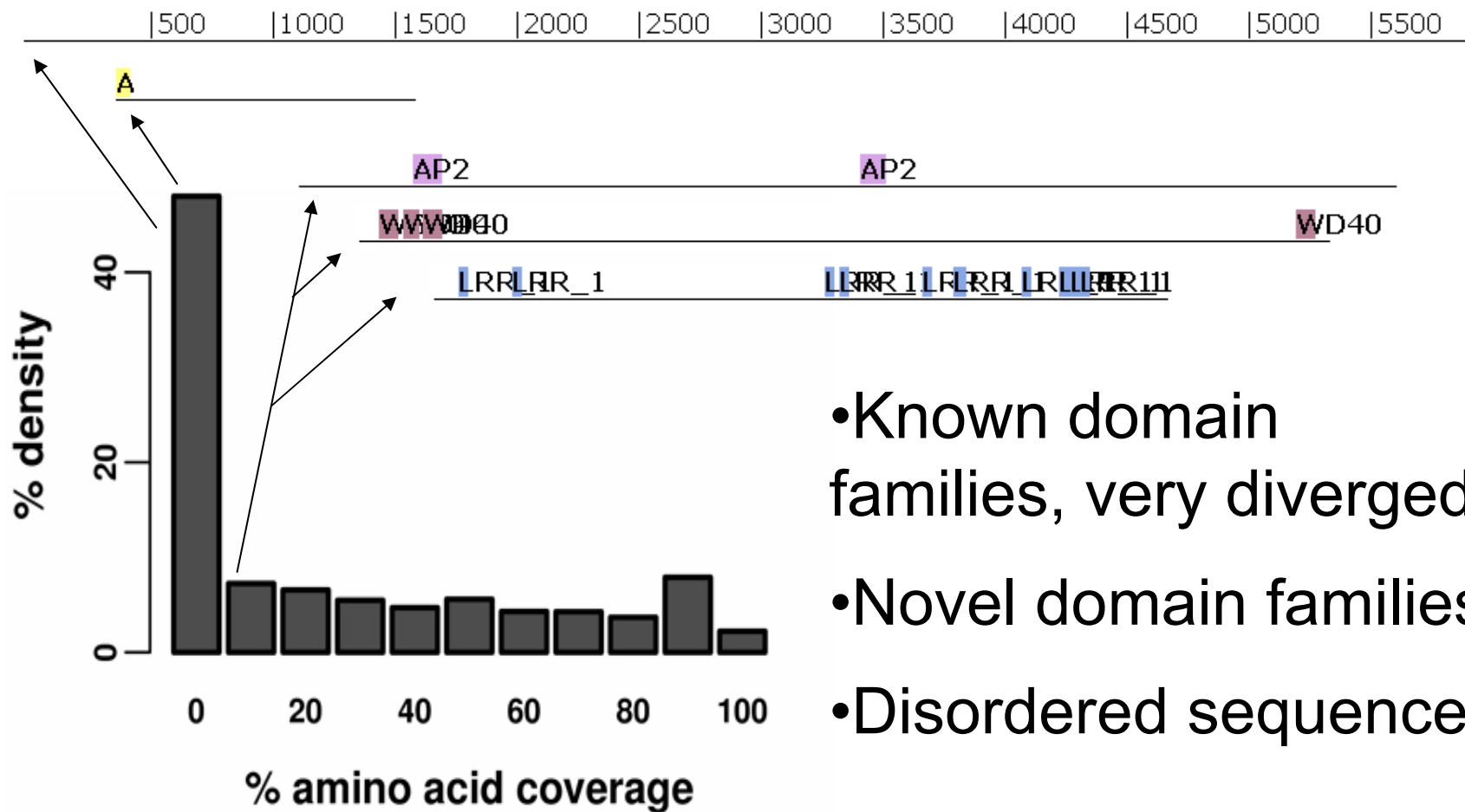
- For new sequences, before experiments start...
- Domains may imply functions
- Experimental alternatives are infeasible as protein databases increase exponentially

# *Plasmodium falciparum*

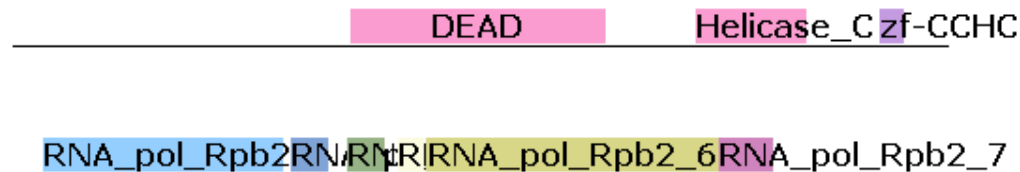
- Malaria
- Diverged eukaryote
  - 80% AT-bias
  - Low-complexity regions
  - Non-photosynthetic plastid
- Annotation
  - 5.5K proteins
  - 45% unknown function
    - 20% unknown in yeast
  - 88% of annotations are bioinformatical



# Poor domain coverage of *Plasmodium falciparum*

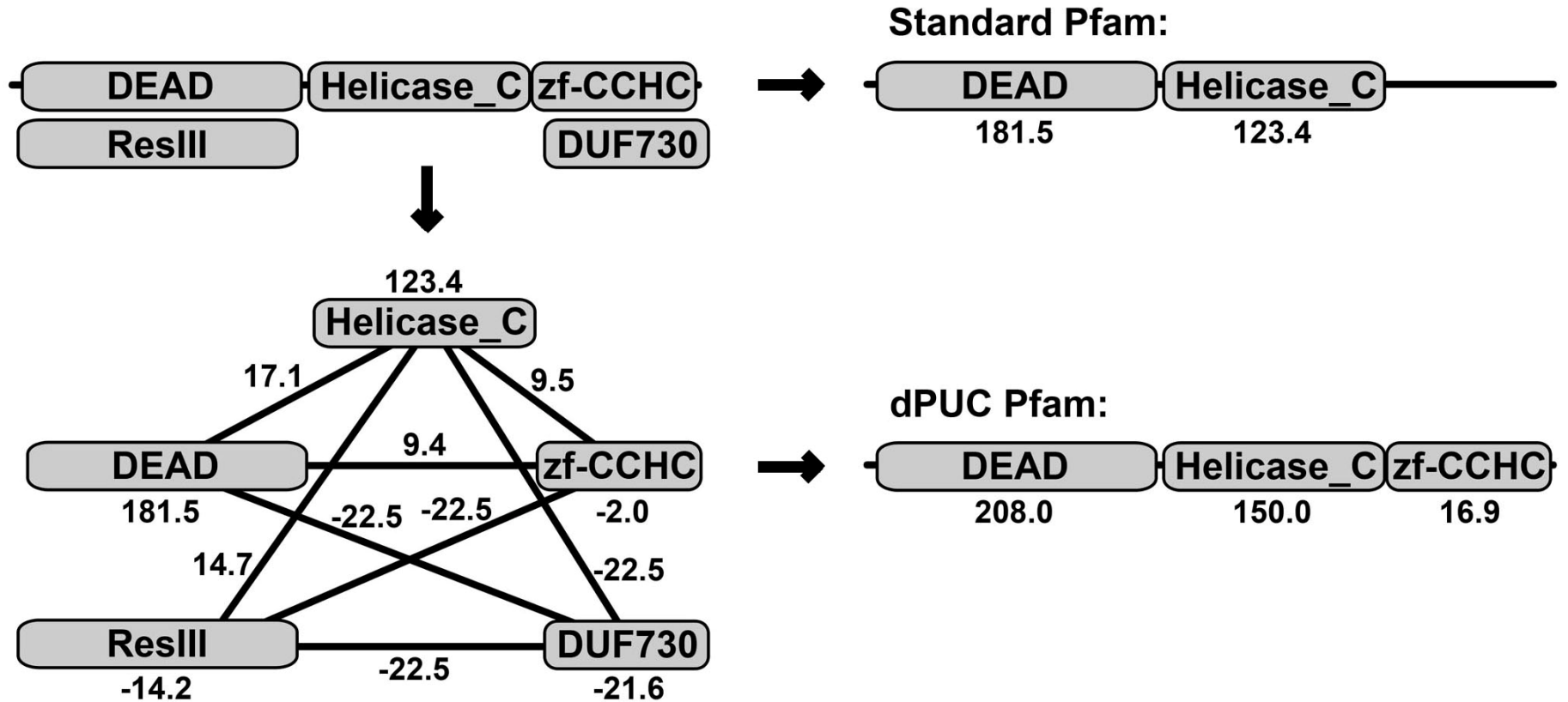


# Domain Prediction Using Context: dPUC



- 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



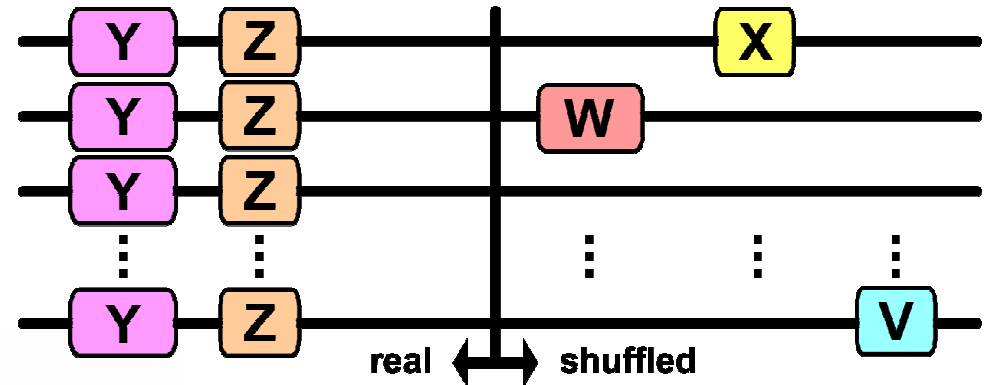


# Improved signal to noise

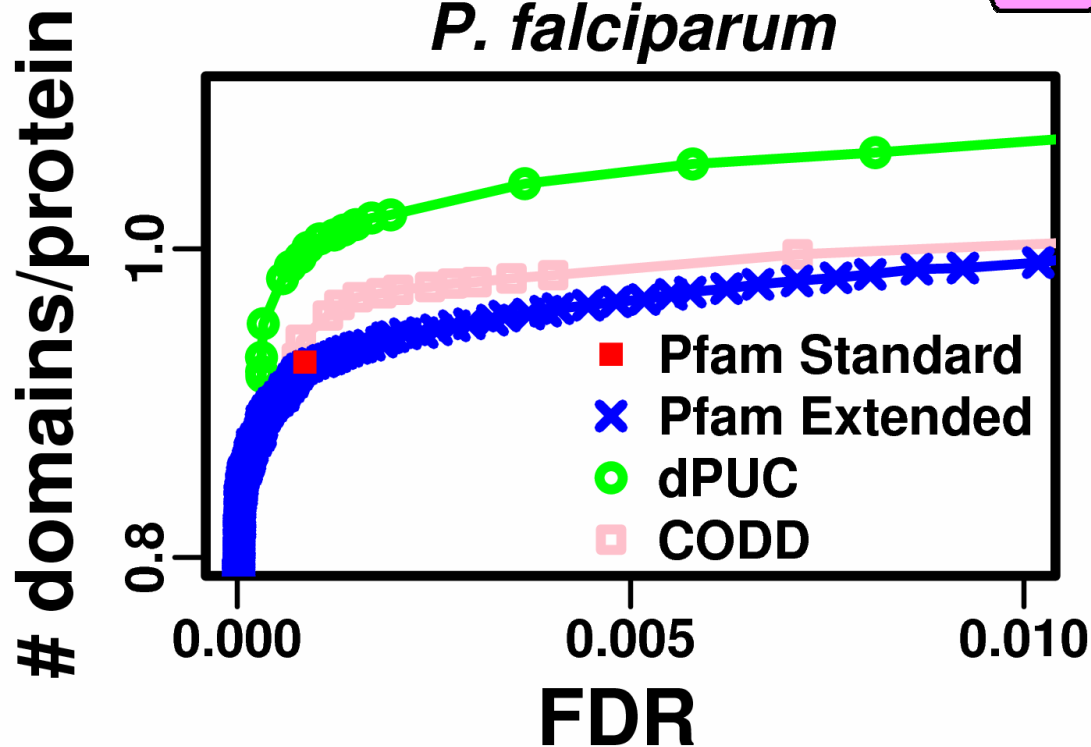
Real protein



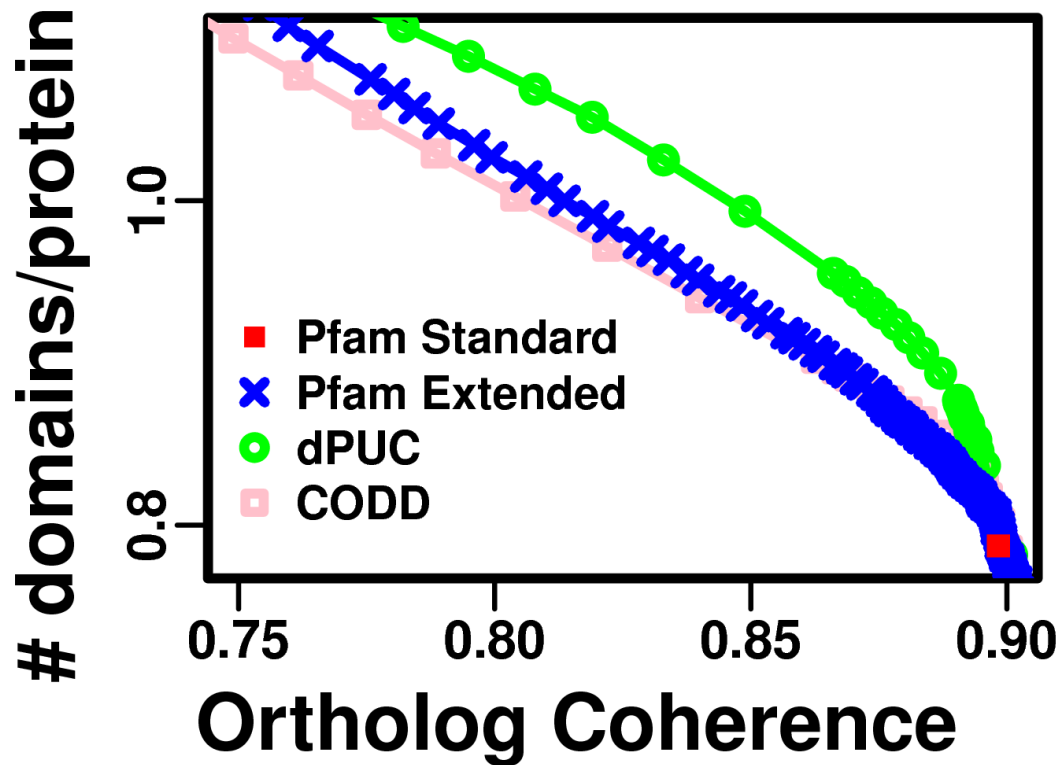
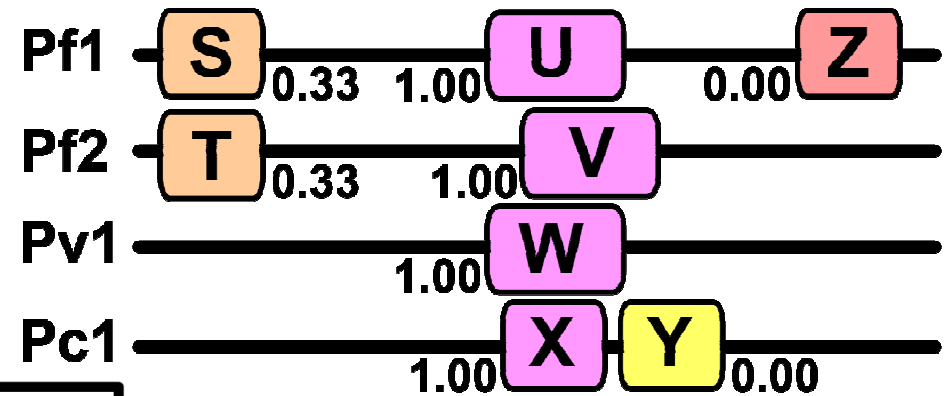
Real protein with shuffled sequences



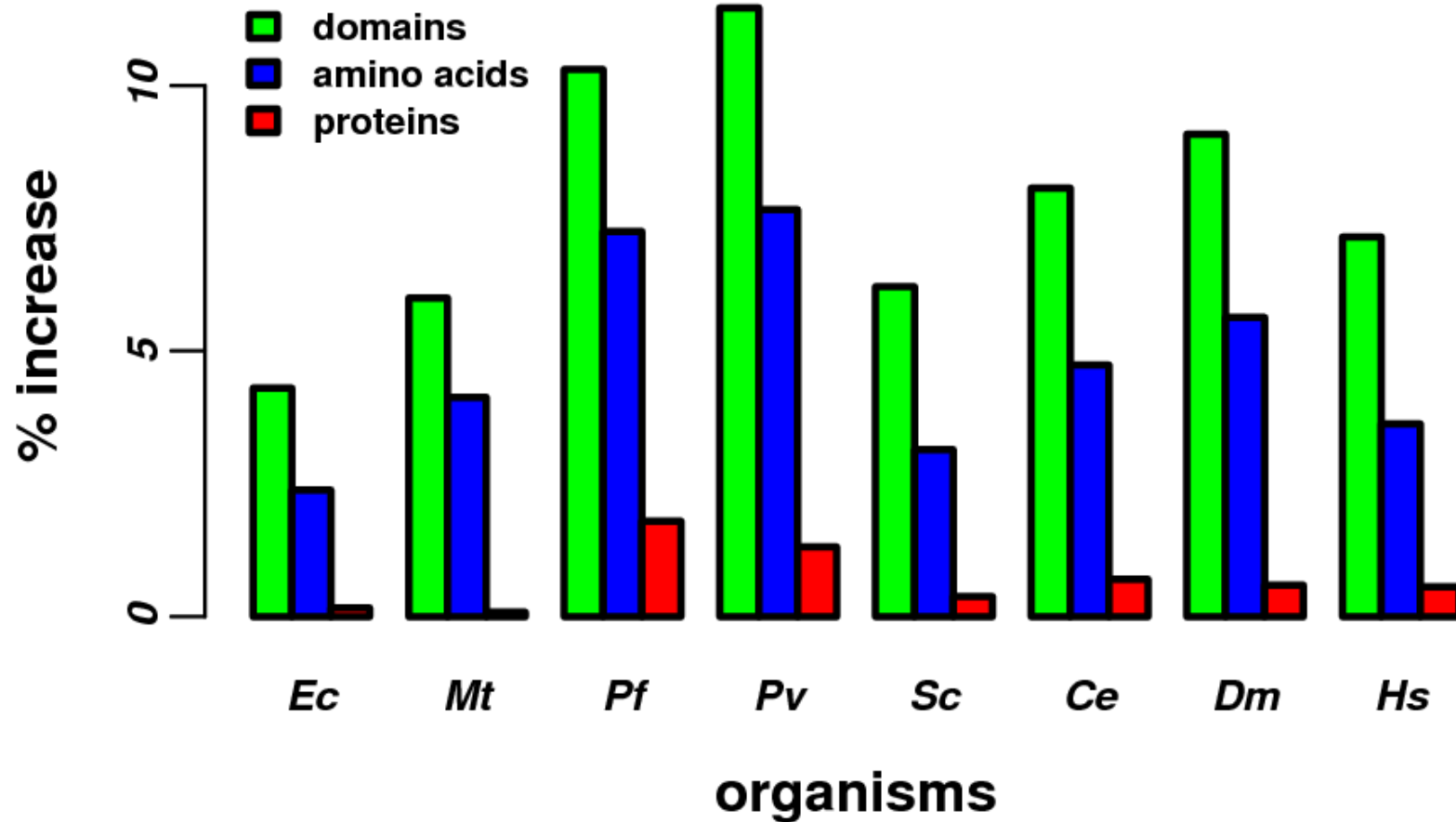
*P. falciparum*



# Improved ortholog coherence on *Plasmodium* species



# dPUC increases coverage



# New predictions

- Phosphatase -> RNA lariat debranching enzyme
- *P. falciparum*

**Standard Pfam**

**dPUC Pfam**

Metallophos

Metallophos

DBR1

- *S. cerevisiae*

**Standard Pfam**

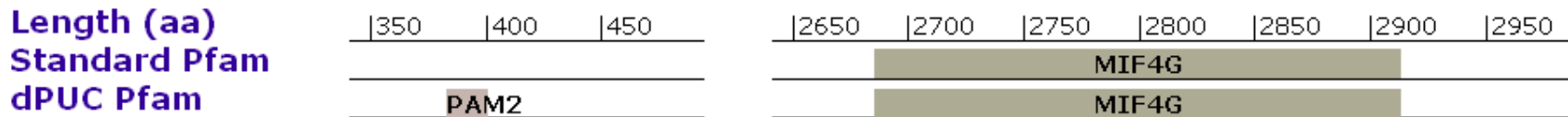
**dPUC Pfam**

MetallophosDBR1

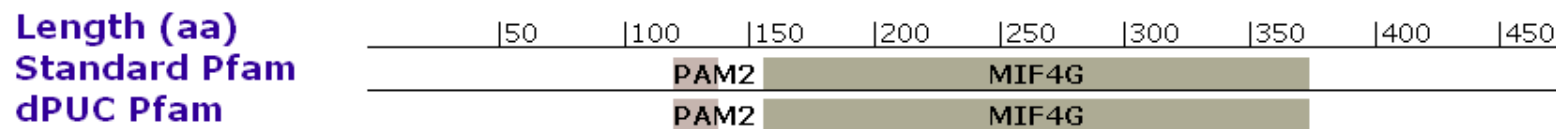
MetallophosDBR1

# New predictions

- MIF4G domain-containing protein -> Poly-A binding protein-interacting protein 1
- *P. falciparum*









- *H. sapiens*

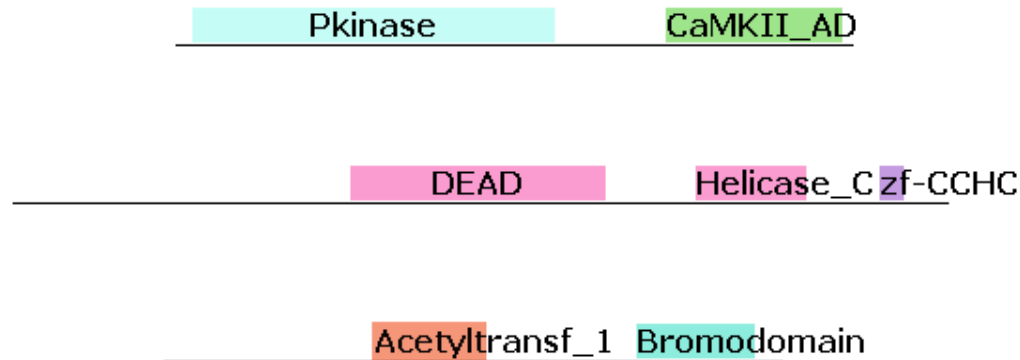


# New predictions

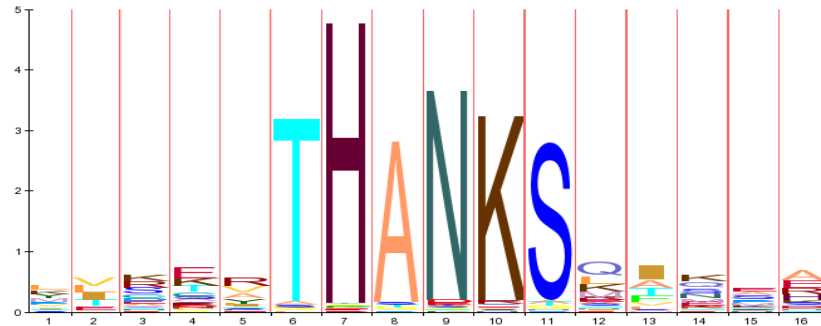
- RNA helicase -> Post-translational mRNA regulation

<b>Description</b>	RNA helicase-1
<b>Organism</b>	<i>P. falciparum</i>
<b>Standard Pfam</b>	
<b>dPUC Pfam</b>	
<b>Description</b>	DDX41_DROME ATP-dependent RNA helicase abstrakt
<b>Organism</b>	<i>D. melanogaster</i>
<b>Standard Pfam</b>	
<b>dPUC Pfam</b>	
<b>Description</b>	DDX41_HUMAN Probable ATP-dependent RNA helicase DDX41
<b>Organism</b>	<i>H. sapiens</i>
<b>Standard Pfam</b>	
<b>dPUC Pfam</b>	

# Domain context



- Complements sequence evidence
- Improves domain predictions
- Works best on *Plasmodium*

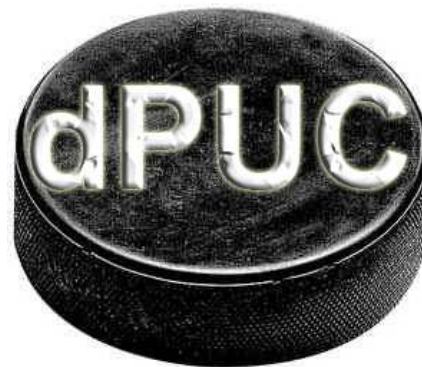


- **Mona Singh & Lab**

- Jesse Farnham
- Dario Gherzi
- Peng Jiang
- Zia Khan
- Anton Persikov
- Jimin Song
- Tao Yue

- **Thesis Committee**

- Leonid Krugliak
- Saeed Tavazoie



[dpuc.princeton.edu](http://dpuc.princeton.edu)

- **Manuel Llinás & Lab**

- Lindsey Altenhofen
- Tracey Campbell
- Erandi De Silva
- Björn Kafsack
- Ian Lewis
- Yael Marshall
- Jessica O'Hara
- Kellen Olszewski
- Heather Painter
- Yoanna Pumpalova

- **NSF GRFP**