

# Functional Annotation

# Functional vs Structural Annotation

- Structure:
  - Location of exons/introns/UTRs
  - Location of regulatory elements
  - Location of repeats
  - Identification of genomic elements by position
- Function:
  - What does it do?
  - Biological implications
  - Genes - what proteins and pathways?
  - Regulatory elements – what do they regulate?

# Outline

- Hidden Markov Models
- HMMER search of Pfam
- Infernal search of Rfam
- Gene Ontology

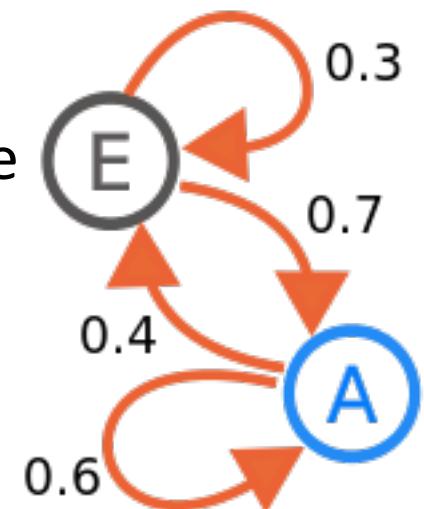
# Pfam and HMMER

# Pfam + HMMER

- There is more to the world than just BLAST (ie traditional sequence alignment)
- The second most popular algorithm is HMMER.
- HMM = Hidden Markov Model
- But to understand that we need to talk about...

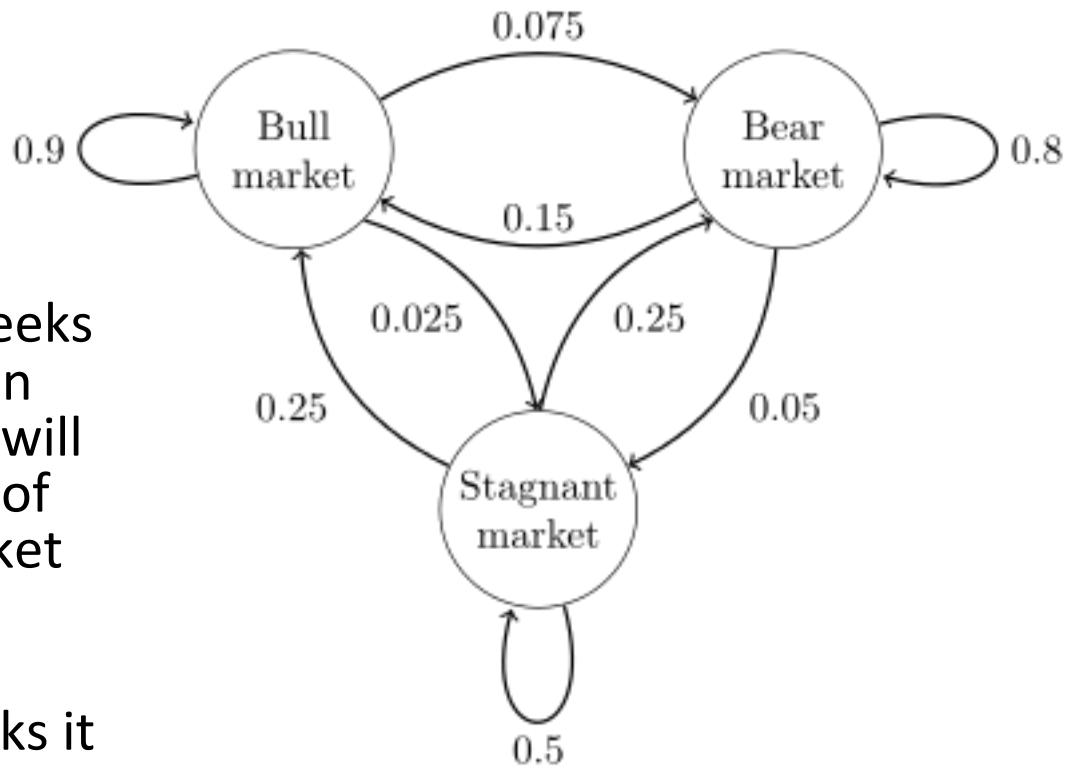
# Markov Chain

- A Markov chain is a random process that undergoes transitions from one state to another on a state space
- Has the property of “memorylessness”
- the probability distribution of the next state depends only on the current state and not on the sequence of events that preceded it
- Called the Markov property
- A Markov chain is a type of Markov Model that is fully observable – we know all the states and probabilities for moving between states



# Markov Chain

- How is it used statistically?
- Possible to calculate:
- the long-term fraction of weeks during which the market is in each state (62.5% of weeks will be in a bull market, 31.25% of weeks will be in a bear market and 6.25% of weeks will be stagnant)
- the average number of weeks it will take to go from a stagnant to a bull market

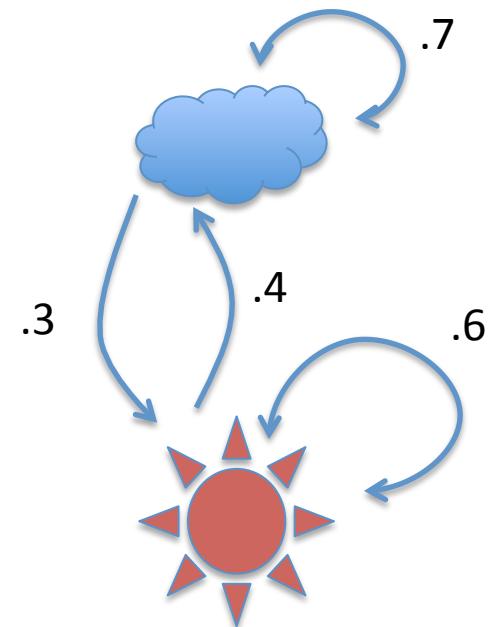


# Hidden Markov model

- The markov chain is only one type of markov model. Another is the hidden Markov model.
- Similar to a Markov chain
- Hidden (unobservable) states
- Example

# Hidden Markov model

- Bob's city weather
- If the day is Rainy
  - 70% chance the next day will be rainy
  - 30% chance the next day will be sunny
- If the day is Sunny
  - 40% chance the next day will be rainy
  - 60% chance the next day will be sunny
- 



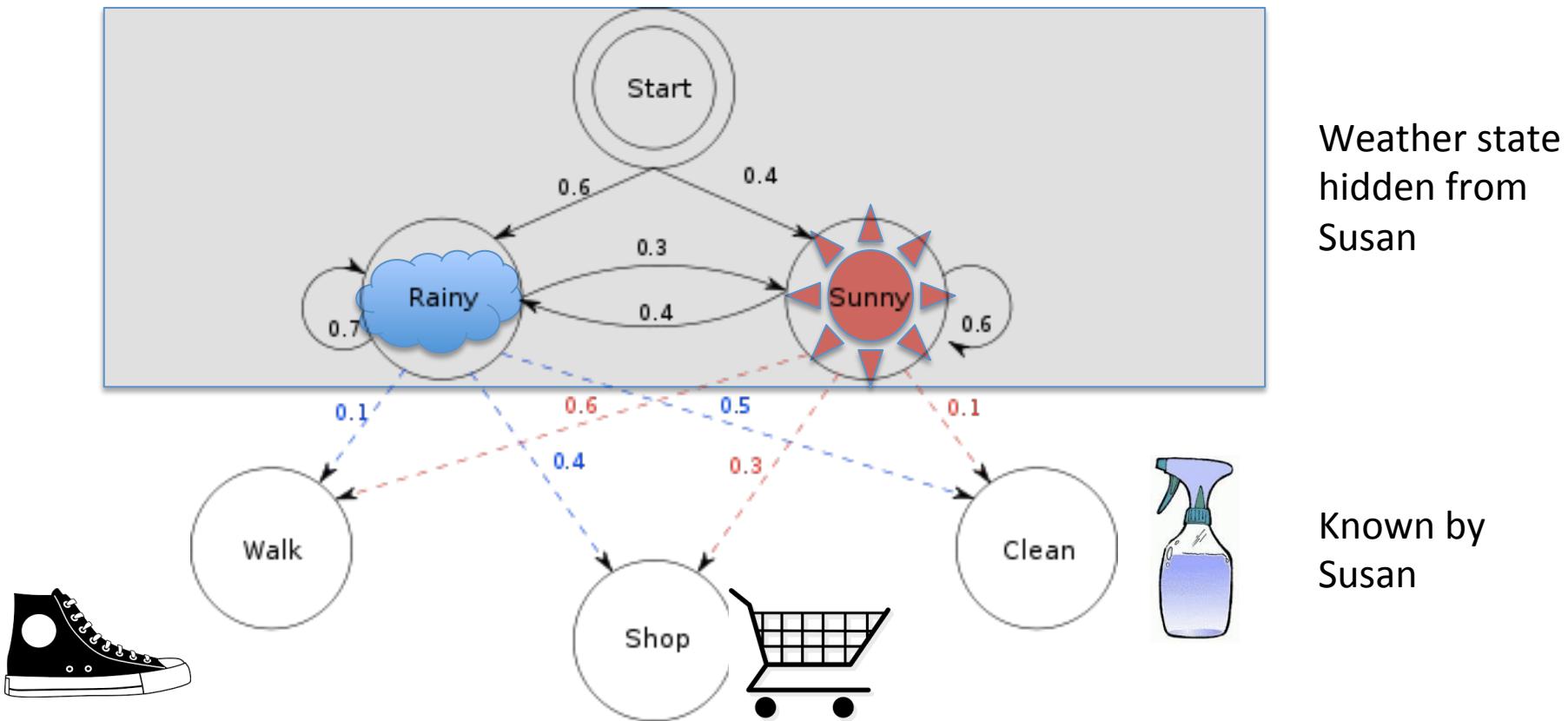
# HIDDEN MARKOV MODEL

- Bob decides what to do based on the weather. He only has three activities: walking, shopping and cleaning.
- If its rainy
  - 10% chance he will walk
  - 40% chance he will shop
  - 50% chance he will clean
- If its sunny
  - 60% chance he will walk
  - 30% chance he will shop
  - 10% chance he will clean



# Hidden Markov Model

- Bob has a friend Susan. Everyday he posts on Facebook weather he is walking, shopping or cleaning. Susan is a mathematician and recognizes this as an HMM.



# HIDDEN MARKOV MODEL – The 3 problems

Using a string of observations about Bob's behavior, what can Susan deduce?

- {walk, clean, clean, shop, walk, walk}
1. Given the model and observations, what is the likelihood of this sequence of observations occurring?
  2. Given the model and observations, what is the likeliest hidden state sequence to have produced this sequence of observations?
  3. Given the observations and the number of hidden states (but not knowing the model), what is the most likely model? (ie Training a model to best fit the observed data.)

# What does this have to do with biology?

- Allow you to incorporate heterogenous types of information for a problem
  - Allow you to add new information more easily.
  - Gene finding. We should account for:
    - splice-site consensus
    - codon bias
    - exon/ intron length preferences
    - open reading frame analysis
  - HMMs provide a conceptual toolkit for building complex models.
- 
- How should the parameters be set?  
How do we weight them?  
How to score?  
How confident that an answer is correct?

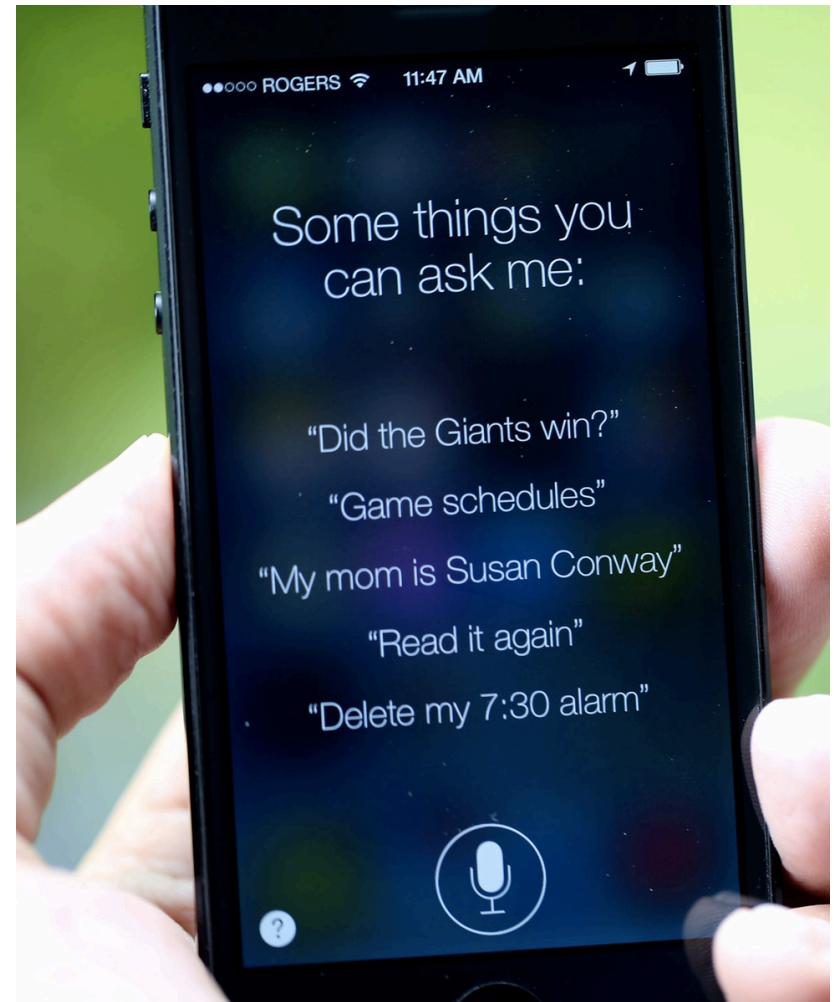
# What does this have to do with biology?

Problems often addressed with HMMs:

- Finding a gene
- Searching for a sequence profile
- Multiple sequence alignment
- Regulatory site identification

Outside of biology, best known for temporal pattern recognition:

- Speech
- Handwriting
- Gesture



# HMM – 5' splice example

- We have a sequence.

Definitely Exon



CTTAGATCGAAATTGATTTCGTAAAACGTTCCCCGG

?????????

Definitely Intron

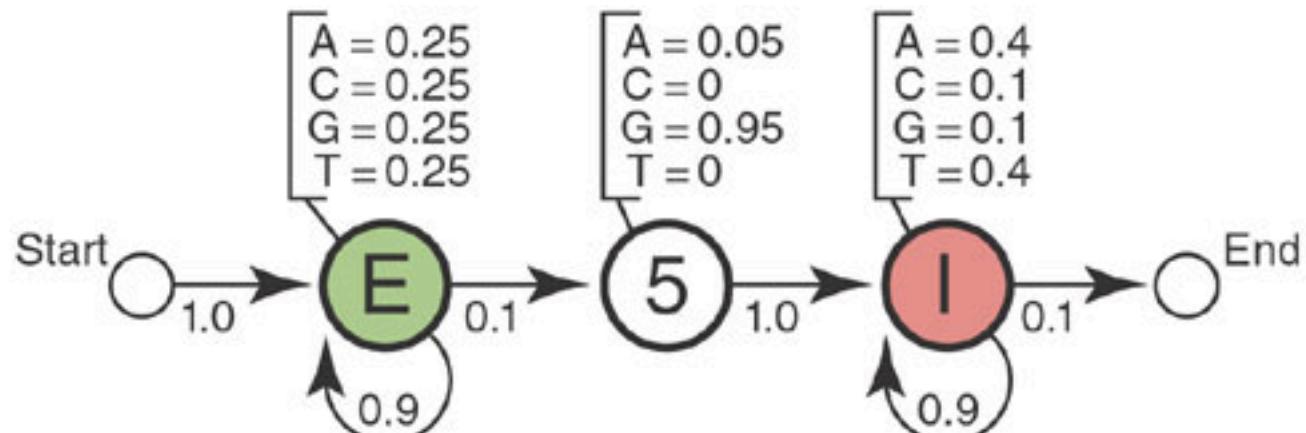


Where is the splice site?

# HMM – 5' splice example

- Lets say we know some information about splicing that will be helpful
- exons have a uniform base composition on average (25% each base), while introns are A/T rich (say, 40% each for A/T, 10% each for C/G),
- the 5'SS consensus nucleotide is almost always a G (say, 95% G and 5% A).
- We can make an HMM.
- We have hidden states: each base is an Exon(E), an Intron(I) or a 5'SS(S)
- We need to find the most likely state that produced the observed sequence

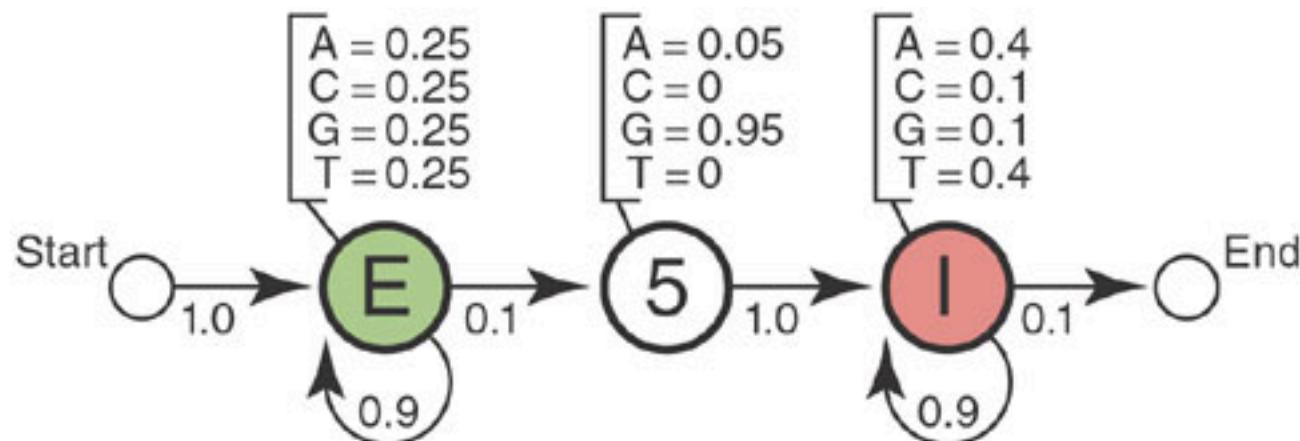
# HMM – 5' splice example



Sequence: C T T C A T G T G A A A G C A G A C G T A A G T C A

Lets test different underlying states to see which is the most likely.

# HMM – 5' splice example



Sequence: C T T C A T G T G A A A G C A G A C G T A A G T C A

State path: E 5 I I I I I I I I log P  
-41.22

Parsing:  
Eddy  
What is a  
hidden  
markov  
model?  
2004

	-41.22	-43.90	-43.45	-43.94	-42.58	-41.71
Eddy	—	—	—	—	—	—
What is a	—	—	—	—	—	—
hidden	—	—	—	—	—	—
markov	—	—	—	—	—	—
model?	—	—	—	—	—	—
2004	—	—	—	—	—	—



- Start with a multiple sequence alignment
- Feed into **hmmbuild**
  - Generate an **hmm profile**
- Calibrate the model with **hmmc\_calibrate**
  - Increase sensitivity of searching
- Search for new homologs that belong to your group with **hmmsearch**



- Why not use BLAST?
- Has much more power in the case of many sequences from the same family – can build a more accurate model of that family by using information about:
  - how conserved each column of the alignment is
  - which residues are most likely at each position
- With a well described protein family, can detect much more remote evolutionary relationships than BLAST.
- Used to be much slower, with new HMMER3 implementation, now is almost as fast as BLAST
- What sorts of databases can we search with HMMER?



- Within a database of protein sequences, many are members of existing protein families and have similar functions. How to organize this information?
- Need to identify protein clusters and to produce multiple sequence alignments.
- The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs).
- Originally published in 1997
- Pfam-A = manually curated family data
- Pfam-B = computationally generated family data



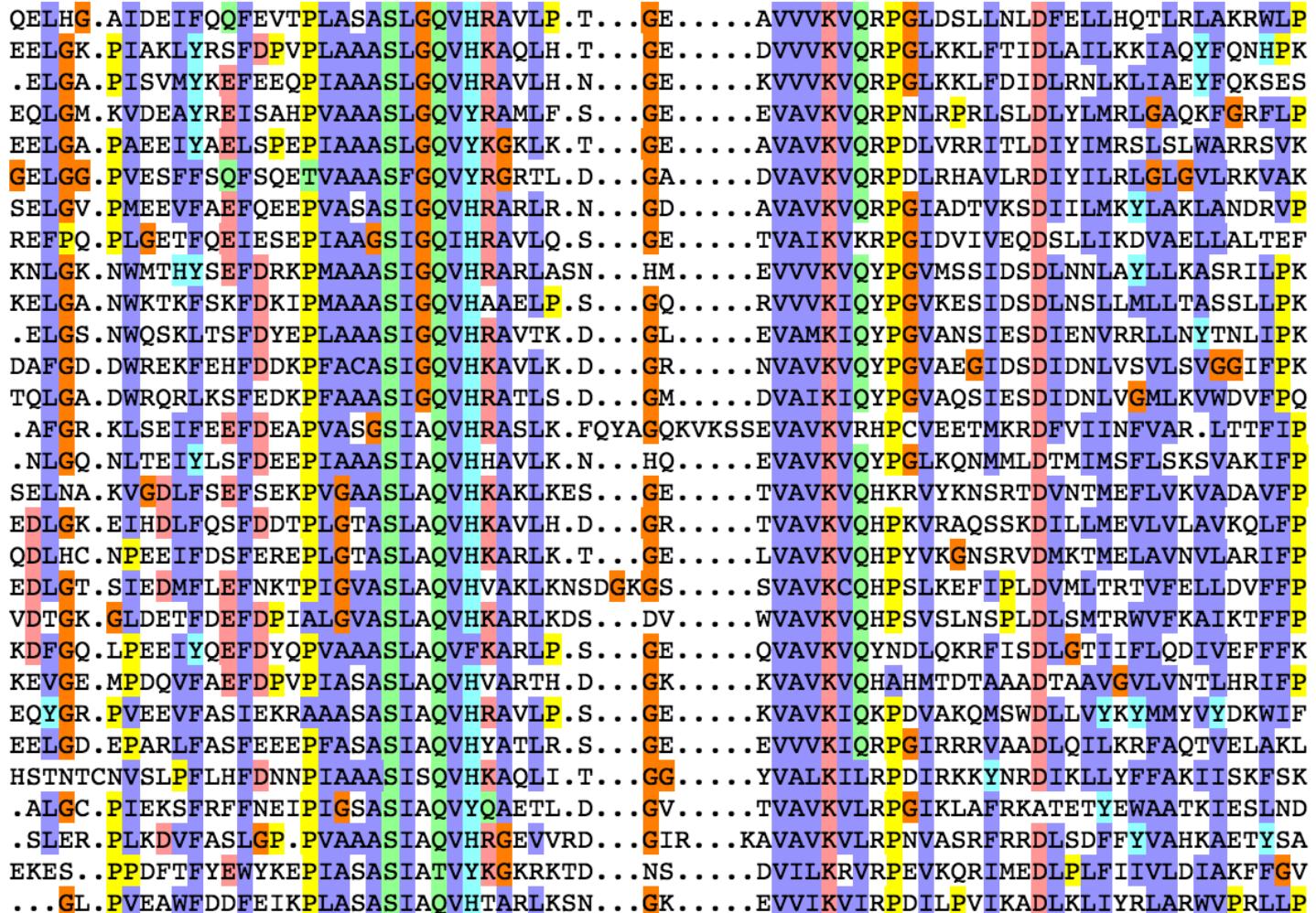
- Currently has 16,306 families (version 30)
- Families are grouped into “clans” - related by similarity of sequence, structure or profile-HMM
- Family information includes gene architecture, structure, sequences from hundreds to thousands of species and interactions.



- HMMER is used by Pfam for two purposes:
  - To construct Pfam clusters (build the model)
  - To detect matches from a given sequence to a cluster (compare a sequence to see how well it fits the model)
- The matching amino acid probabilities from the HMM can be visualized as a logo

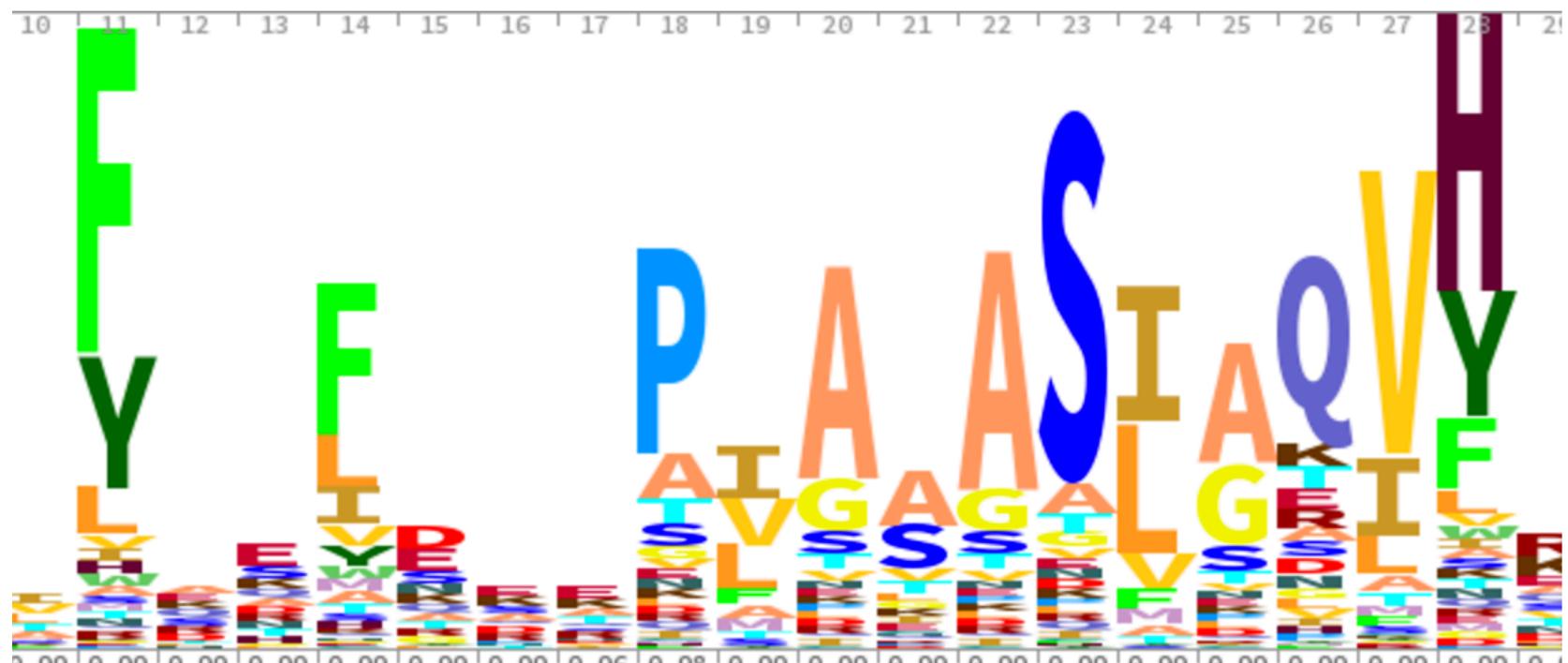
# Multiple Sequence Alignment

Y095\_SYN3/115-238  
Y1770\_SYN3/142-261  
B9DGY1\_ARATH/248-365  
Y1919\_SYN3/127-246  
Y005\_SYN3/161-279  
O80962\_ARATH/256-373  
O27682\_METTH/119-238  
Y889\_SYN3/100-218  
ABCI\_SCHPO/284-401  
COQ8\_YEAST/176-292  
Q9SBB2\_ARATH/284-398  
COQ8\_CAEEL/417-533  
Q9VYI6\_DROME/336-452  
Q3ECK9\_ARATH/268-392  
F4ID59\_ARATH/156-271  
O17735\_CAEEL/142-259  
ADCK1\_HUMAN/143-259  
Q9W133\_DROME/137-253  
MCP2\_YEAST/166-289  
MCP2L\_SCHPO/168-286  
Q9VTG5\_DROME/162-278  
Y2090\_ARATH/147-268  
YF9E\_SCHPO/167-287  
Y647\_MYCTU/150-271  
Q9ZCP5\_RICPR/34-153  
H2VFS0\_ZYMMO/111-228  
Q89WD1\_BRADU/112-231  
Y445\_PBCV1/90-208  
UBIB\_SHIDS/115-232



# Becomes a Profile

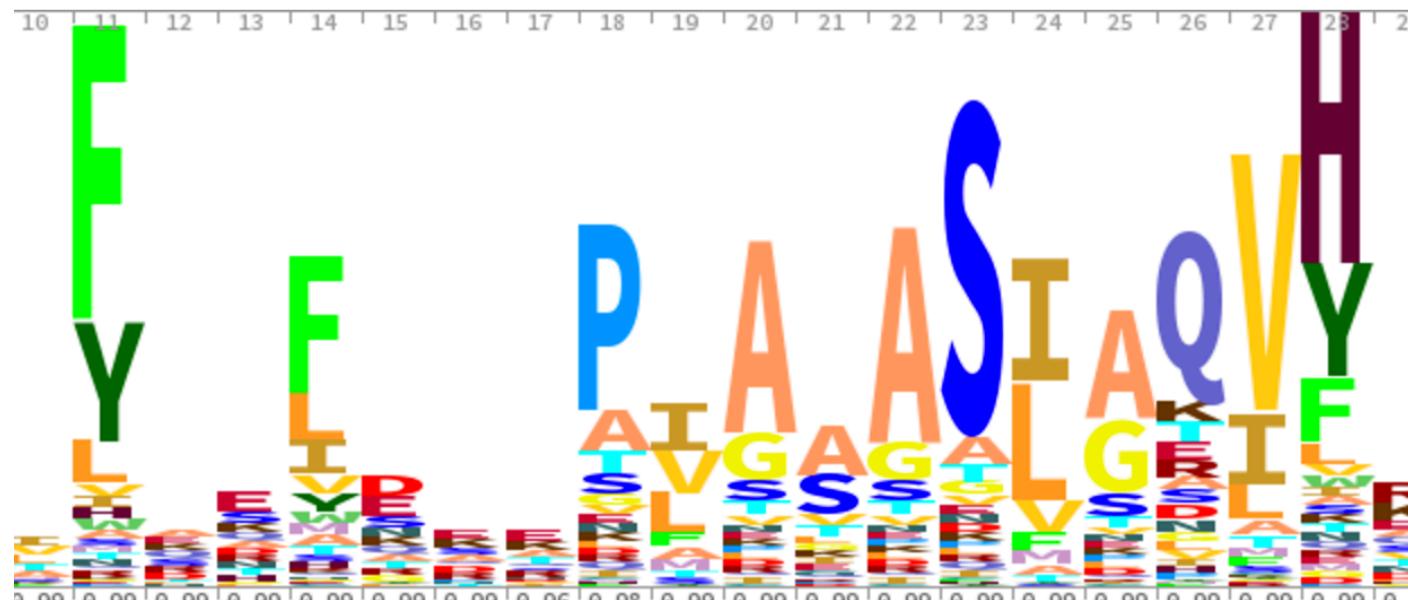
- HMM represented as a logo



Height = letter probability

# Pfam

- Take all 16,306 profiles
- Use hmmsearch to compare your sequences to these families
- E-values –just like BLAST



# Example:

- Search green ash protein

ALCLIMLAHSGGGAAISPNSVNTTRPNLPTINDSKQIENSTTPPPTQDQSYSCVCNKAFASYQALGGHKASHRKNATASDDG  
NHSTSTTTAAASTASNVSAALNPRGRLHECSICHKSFPTGQALGGHKRRHYEGIIGGGSSKSSVTSSDGGASSHAPRDFDLNLPATP  
EFQLELTVDVCVKKSQFVGDQEVESPMPFKKPRT.PT.FGERF

- Results

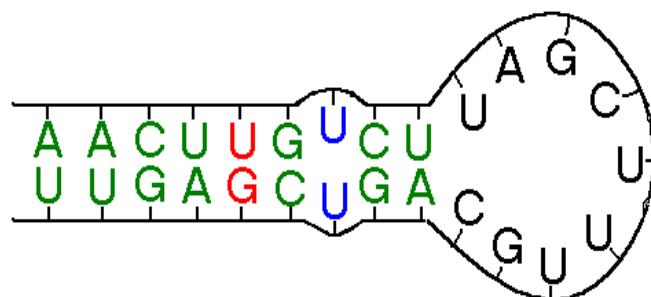
Family	Description
<u>zf-C2H2_6</u>	C2H2-type zinc finger
#HMM	heCdeCsksFpS1qaLggHkksHrk
#MATCH	+ C++C+k F S+qaLggHk+sHrk
#PP	78*****8
#SEQ	YSCSVNCNAFASYQALGGHKASHRK

Entry type	Clan	Envelope		Alignment		HMM		HMM length	Bit score	E-value
		Start	End	Start	End	From	To			
Domain	<u>CL0361</u>	53	79	54	78	2	26	27	45.9	3.2e-12

Infernal + Rfam



- Infernal ("INFERence of RNA ALignment")
- Tool for searching for RNA structure and sequence similarities
- Uses information about sequence AND structure
- Also uses HMM
  
- Start with:
  - Multiple sequence alignment
  - Special annotation of bases that are paired to create the secondary structure





Version 12.1  
2,474 families

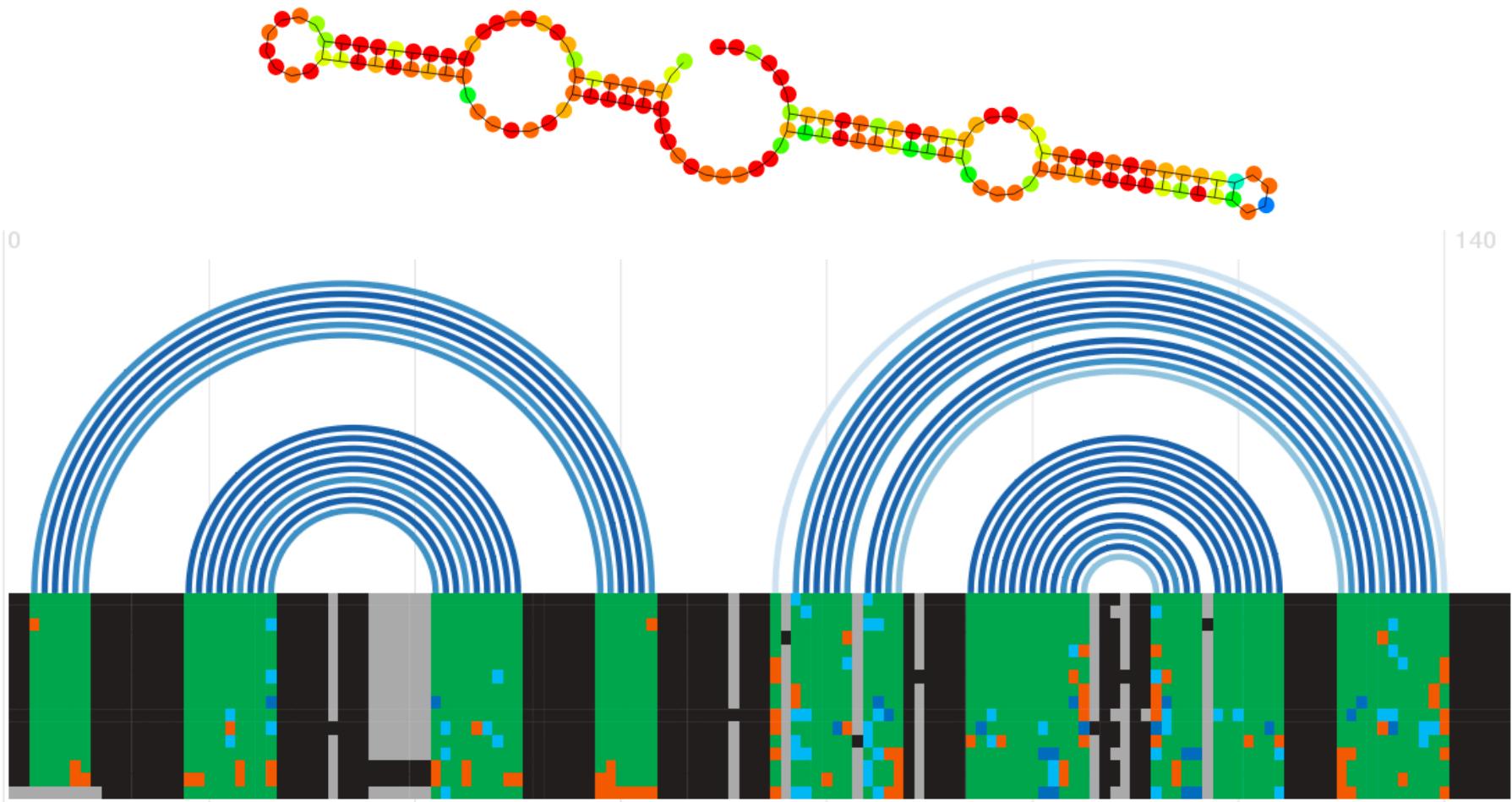
- Contains RNA families with conserved secondary structure
- RNA types:
  - non-coding RNA genes
  - structured cis-regulatory elements
  - self-splicing RNAs
- Grouped into RNA families, each represented by:
  - multiple sequence alignments
  - consensus secondary structures
  - covariance models (CMs).
- Search with Infernal



# Example snoRNA

AAAGCAGGUUGCAAUUACAGUGCUCUCAUUU.GUG.....GAAGUACUGCCAUUAUCCUGCUGAAAGAA.AAGC.CGUGUU.AAUCA.UUUUUGAUUUUGCCUU.UA  
AAAGCAGGUAGCAAUUACAGUGCUCUCAUUU.GUG.....GGAGUACUGCCAUUAUCCUGCUGAAAGAA.AAGC.CAUGUU.GGUUG.UUUCUGAUUUUGCCUU.U-  
-----AGCAAUUACAGUGCUCUCAUUU.GUG.....GGAGUACUGCCAUUAUCCUGCUGAAAGAA.AAGC.CAUGUU.GGUUG.UUUCUGAUUUUGCCUU.U-  
AGAGCAGGUUGCAAUUACAGUGCUCUCAUUU.GUG.....GAAGUACUGCCAUUAUCCUGCUGAAAGAA.AAGC.UAUGUU.GAUCA.UUUUUGAUUUUGCCUU.C-  
AAAGCAGGUUGCAAUUACAGUGCUCUCAUUU.GUG.....GAAGUACUGCCAUUAUCCUGCUGAAAGAA.AAGC.UGUGUU.GAUCG.UUAUUGAUUUUGCCC.UA  
AAAGCAGGUUGCAAUUACAGUGCUCUUCGUUU.GUG.....GAAGUACUGACAUUAUCCUGCUGAAAGAA.AAAC.AGUGUU.GAUCA.UUUUUGAUUUUGCCUC.UC  
AAAGCAGGUUGCAAUUACAGUGCUCUUCUUU.GUG.....GAAGUAUUGACAUUAUCCUGCUGAAAGAA.AAUC.UGUGUU.GAUCGuUUUUUGAUUUUGCCAU.UUa  
UACGCAGGUUGCAAUUACAGUGCUCUUGUUU.GGG.....GAAGUACUGCUGUUAUCCUGCUGAAAGAC.AAGC.UGUGUU.AGUCA.UUUUUGAUUUUGCCUU.UA  
AAAGCAAGCUGCAAUUACAGUGCUCUCAUUU.GUGaaaacUAAAACUGCCAUUAUCCUGCUGAAAGAA.AAGC.UGUGUU.AAUGA.UUUUUGAUUUUGCCUU.UG  
AAAGCAGGCUGCAAUUACAGUACUUCAGUUU.GUG.....GAAGUACUGCCAUUAUCCUGCUGAGAGAAAGC.CAUGUU.GGCCG.GCUCUGGUUUUGCCUC.U-  
UGAGCAGGUUGCAGUCCAGUCUUUGUUUcGUG.....GGAGUGCUGGCAUAACCCUGCUGAAAACA.AAUA.UGUGCC.AAUCA.UUUUUUUAUUUACCUCaUU.  
UAAGCAGGUUGCAAUUACAGUGCUCUCAUUU.GUG.....GAAGUACUGACAUUAUCCUGCUGAAAGAA.AAUCAUGUGUG.GAUCA.UUUUUGAUUUUGCCUU.UG  
AAAGCAGGUUGCAAUUACAGUACUUCAUUCU.GUG.....GAAGUAUUGCCAUUAUCCUGCUGAAAGAA.AAGC.CGUGUUuAAUCA.UUUCGGGUUUUGCCUG.UA  
AAAGCAGGUUGCAAUUACAGUGCUCUCAUUU.GUG.....GAAGUACUGACAUUAACCCUGCUGAAAGAA.AAUG.UGUGUC.GAUCA.UUUUUGAUUUUGCCUU.UA  
AAAGCAGCUGGAAUUGCAGUGCUCUCAUUU.GUGaaaacUAAAACCAUCAUUAUGCUGCUGAAAGAA.AAGC.UGUUUU.AAUGA.UUUUUGAUUUUGCCUU.UG  
AAAGCAGGUUGCAAUUACAGUGCUCUUAUUU.GUG.....AAAGUACUGUCAUUAUCCUGCUGAAAGAA.AAGC.UGUGUU.GGUCC.UUUUUGAUUUUGCCAC.UG

# Example snoRNA



Hidden markov model uses information about sequence conservation and structure conservation

# Ontology

# Ontology

- Roots in philosophy – how we conceptualize and specify knowledge (Aristotle)
  - Very useful for organizing information with computers
  - This is a very big area of thought and utility, we're going to focus on a relatively simple example:
  - Controlled Vocabulary
- Example:
- Wine
    - White Wine
    - Rose Wine
    - Red Wine
      - Beaujolais
      - Red Burgundy
      - Red Zinfandel
      - Merlot
      - Syrah (Shiraz)

# From NCBI SRA for Arabidopsis

I want sequences that relate to flower structures. I have to manually look through a list and select:

- Inflorescense
- Inflorescence
- Immature inflorescence
- Flower
- Flowers
- Pistils pollinated for 8 Hours
- 3xHA\_inflorescence\_biological\_replication1
- 3xHA\_inflorescence\_biological\_replication2
- 3xHA-VvCEB1-OX\_inflorescence\_biological\_replication3

# Plant Structure Ontology

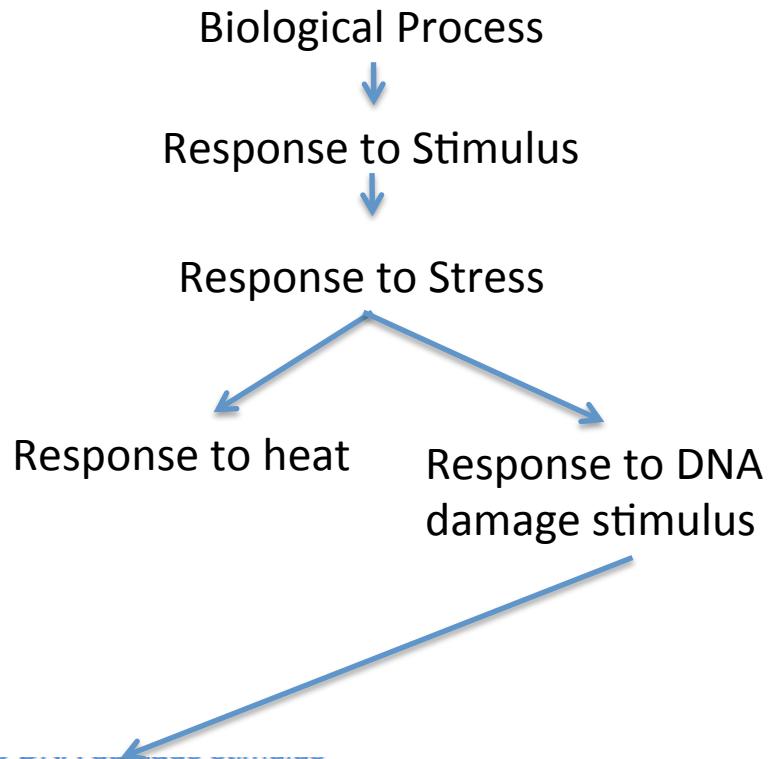
- All these would be coded in a computer readable structure:
  - Inflorescence
  - Flower
    - Gynoecium (Pistil)
    - Androecium
    - Perianth



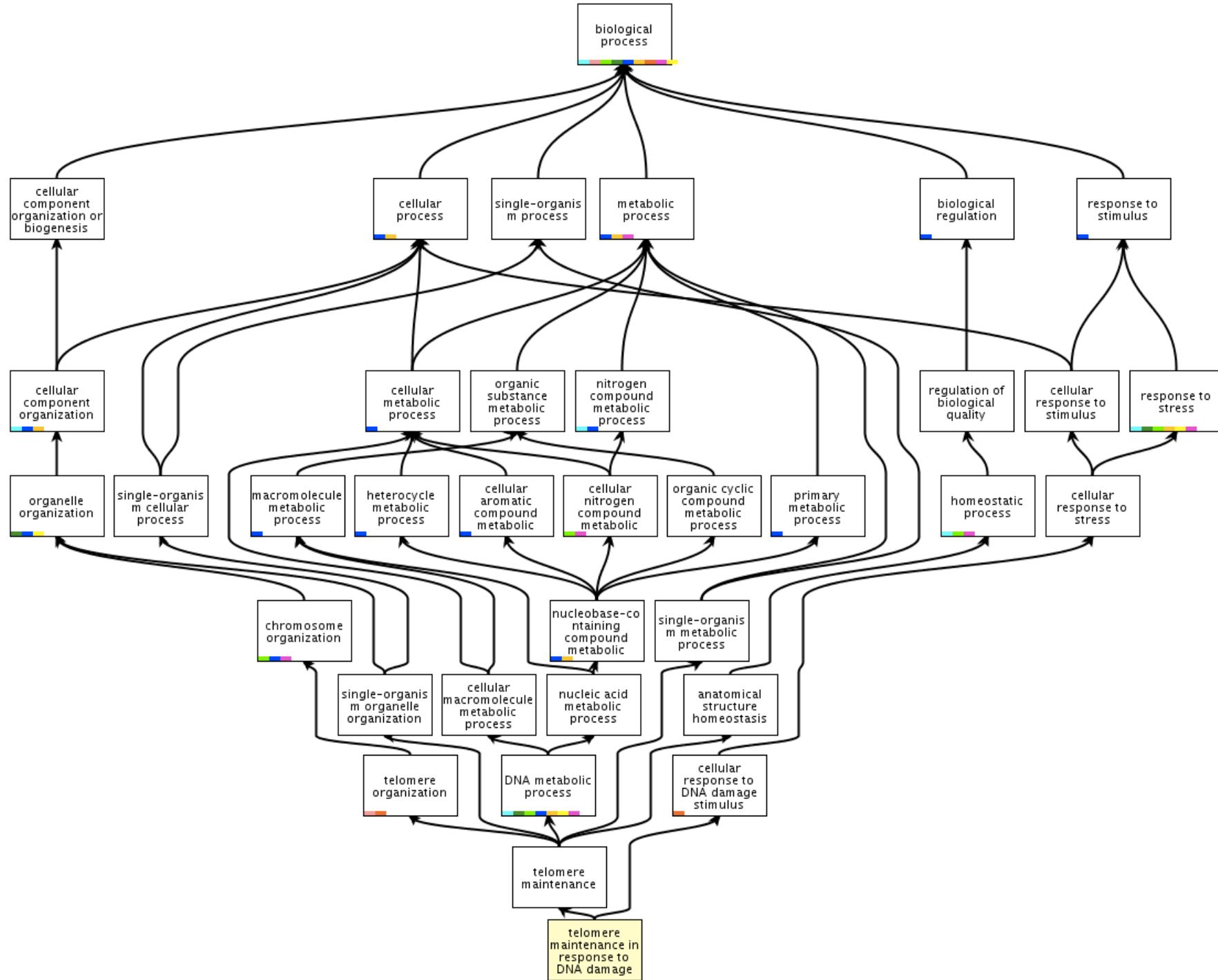
# GENEONTOLOGY

Unifying Biology

- Used for annotating genes
- Three sections
  - Biological Processes
  - Metabolic Functions
  - Cellular Components
- Each section is formed as graph or network of terms



GO:2001022 positive regulation of response to DNA damage stimulus  
GO:2001020 regulation of response to DNA damage stimulus  
GO:1990248 regulation of transcription from RNA polymerase II promoter in response to DNA damage  
GO:0031297 replication fork processing  
GO:0042770 signal transduction in response to DNA damage  
GO:0043247 telomere maintenance in response to DNA damage



# Uses of Gene Ontology

- Example (Without an ontology)
  - You search for genes that are induced by stress
  - The computer searches through a database of genes
  - Some are annotated “response to drought”
  - The computer does not return them because it does not see the word “stress” anywhere
- Example (With an ontology)
  - You search for genes that are induced by stress
  - The computer searches through a database of genes
  - Some are annotated “response to drought”
  - The computer looks that up in our vocabulary and figures out that drought is a type of abiotic stress, which is a type of stress
  - The computer now returns those genes to you
- Knowledge has become more accessible!!!

# Uses of Gene Ontology

- Finding all members of the same biological process or pathway
- Finding high level patterns of metabolic or biological activities
- Looking for statistical enrichment of GO terms
  - E.g. From the control to the treatment, the occurrence of lignan production related genes increases
- Tools:
  - GO home page
  - BinGO cytoscape plugin

# Example

