

# Genome-Wide Compilation of *Arabidopsis* Polyadenylation Signals Using Cray SV1 Accelerated Pattern Recognition

---

Quinn Li (Miami University)

Dave Strenski (Cray)

Eric Stahlberg (OSC)

Johnny Loke (Miami University)

Chris Wood (Miami University)

# *What is Arabidopsis ?*

**Mustard**



**Cauliflower**



**Cabbage**



**Canola**



*Arabidopsis*

(mouse ear cress)

# Why Arabidopsis?

---

- Model organism
- Simple, but has everything that makes up a plant
- Small genome and fast growing, which is good for genetic studies
- Genome sequenced

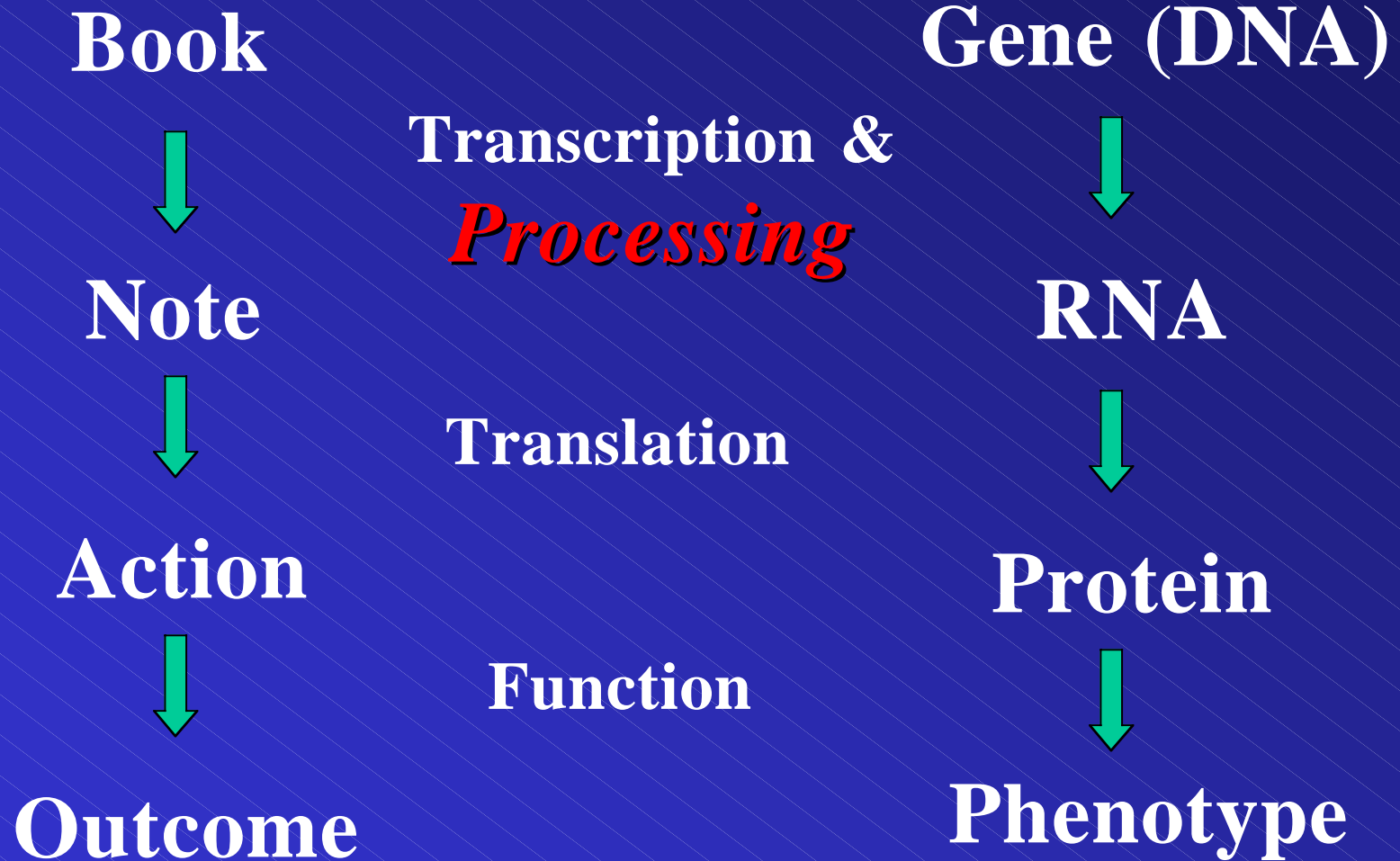
# What is Genome?

---

- The total genetic content of an organism
- Make up of DNA
- DNA use only four characters: A, T, G, C
- Different orders, “sequences”, of the characters become genes, “paragraphs”
- The “words” in the paragraph are the special domains defining the cell’s function

# Gene Expression

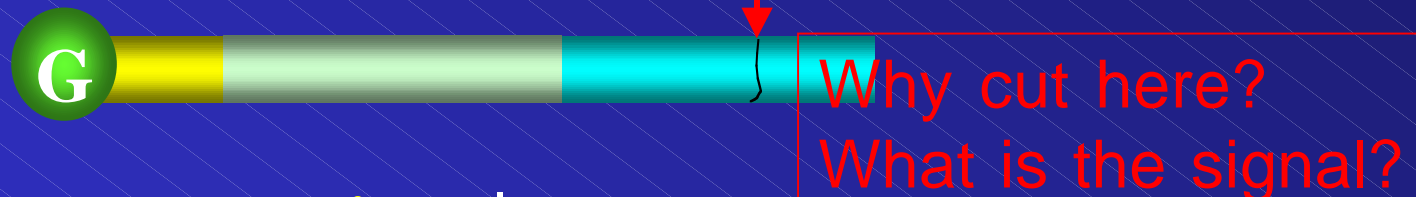
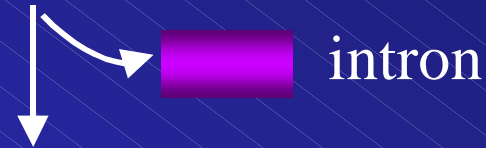
---



# RNA processing



Capping / splicing



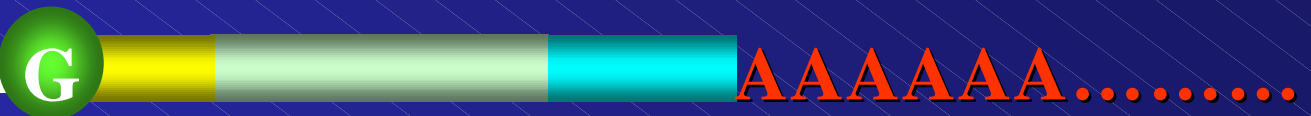
Polyadenylation

mature mRNA



# Database Available (cDNA library)

---

mature mRNA 

Reverse transcription



cDNA

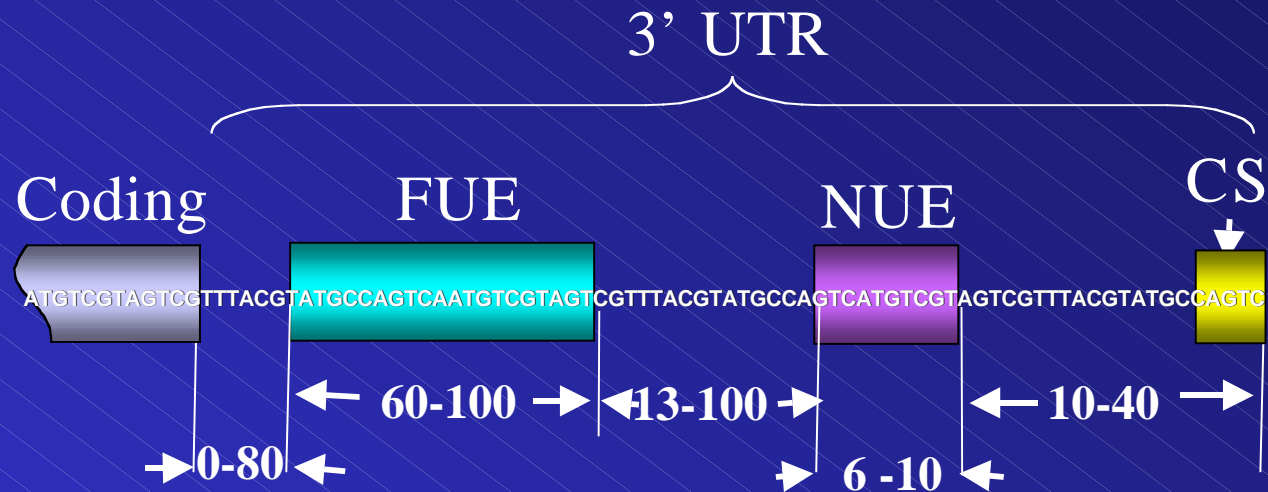


Poly(A) signals

Total 16,225 cDNA in the database



# Working model of poly(A) signals



**NUE: Near Upstream Element**

**FUE: Far Upstream Element**

**CS: Cleavage Site**

**3'UTR: 3' UnTranslated Region**

**Coding: code for protein**



# Current Project

---

- Find the consensus sequences that act as signals
- This happens to most of the 26,000 genes (mRNA) in Arabidopsis/all plants
- The signal should have some commonality: pattern or similarity?

# Current Approaches

---

- Select sequences
- Align sequences (exponential problem with number of sequences)
- Create Hidden Markov Models and find consensus
- **PROBLEM: too many sequences to align!**
- **PROBLEM: need to know what to look for first**

# Using Cray Bioinformatics Libraries to Meet Challenges

- Why Cray libraries
  - Fast and convenient programming interface
  - Large, shared memory, address space for programs
  - Large secondary storage space can store several sets of information at one time

# What are Cray Bioinformatics Libraries?

- A set of routines that accesses the special bit manipulation hardware: popcnt, leadz, vector shift, bit matrix multiply (BMM)
- Rapidly search for sequences (patterns), reverse sequences, and slightly off sequences
- Works for DNA and proteins
- Defined API, open source

# Approach to Locating Poly(A) Signals

We do not know the pattern nor location

Our problem has 16255 sequences ranging in length from 50 to 3118 characters, with the average being 231, total characters about 3.7 million

# Approach to Locating Poly(A) Signals

Here's an example: (next set of slides)

Each line is a different sequence

**Green** is the NUE region of interest

**Red** is a matching pattern

Sequences all end with cleavage site

AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
CGTAATACAAAATAGTACCAAAT  
GTATACATAAATAGACAAGAT  
AATACAGATTCAGAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC



**AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
CGTAATACAAAATAGTACCAAAT  
GTATACATAAATAGACAAGAT  
AATACAGATTCAGAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC**

AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
CGTAATACAAAATAGTACCAAAT  
GTATACATAAATAGACAAGAT  
AATACAGATTCAAGAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC

AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
CGTAATACAAAATAGTACCAAAT  
GTATACATAAATAGACAAGAT  
AATACAGATTCAAGAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC

< AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
<< ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
> CGTAATACAAAATAGTACCAA  
GTATACATAAATAGACAAGAT  
AATACAGATTCAAGAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC

AGTCGTCGTTCTCTAAATAGGCTAGC  
AGCTTCGCTATCTAGCTAGCTAGTACCC  
ATGACATATAACAGAAATAGCATAT  
GTTACATGTACATATCTATAAGATACTAC  
CGTAATACA AAATAGTACCAA  
GTATACATAAATAGACAAGAT  
AATACAGATTCA GAGACATATACACAG  
ATACAGATTCCACAGAAGATACAGAT  
ATACGATCATGCTATACATATCGATC

# Approach to Locating Poly(A) Signals

Looking at all alignments, what pattern is the most common in the NUE region

# Read and Compress the Data

The FASTA file has the following format:

> Sequence One Title

```
AGTCTCTGATGACTGATCATGATC  
TCAGACGCTACCGTACGACTCTAC  
CGTAGATC
```

> Sequence Two Title

```
AGCATCATCAGCTTATGACGGCTA  
ACGTATTGATTCAATCTAC
```



# Read and Compress the Data

Use `cb_read_fasta` to read the file

Returns data array with all the sequence information stored with sequences starting on multiples of 4 word boundaries

Returns pointers to a structure holding information about the start and length of the sequences, and header information

# Read and Compress the Data

A = hex:41 = binary:0100 0001

C = hex:43 = binary:0100 0011

G = hex:47 = binary:0100 0111

T = hex:54 = binary:0101 0100

a = hex:61 = binary:0110 0001

c = hex:63 = binary:0110 0011

g = hex:67 = binary:0110 0111

t = hex:74 = binary:0111 0100

# Read and Compress the Data

ACGT

0100 0001 0100 0011 0100 0111 0101 0100

ACGT = 00 01 11 10

Use `cb_read_fasta` and `cb_compress`

One forth the storage

One forth the I/O to move around

# Find the most Common Pattern

**ascii          binary          decimal**

**AAAAAA = 000000000000 = 0**

**AAAAAC = 000000000001 = 1**

**AAAAAT = 000000000010 = 2**

**GCTAGC = 110110001101 = 3469**

**GGGGGC = 111111111101 = 4093**

**GGGGGT = 111111111110 = 4094**

**GGGGGG = 111111111111 = 4095**

# Find the most Common Pattern

**AAAAAA = pat\_freq( 0) = 0**

**AAAAAC = pat\_freq( 1) = 0**

**: : :**

**GCTAGC = pat\_freq(3469) = 0**

**: : :**

**GGGGGT = pat\_freq(4094) = 0**

**GGGGGG = pat\_freq(4095) = 0**

# Find the most Common Pattern

AGTCGTCGTTCTCTAAATAG**GCTAGC**

AAATAG=0    AATAGG=0    ATAGGC=0

AGTCGT=0    AGGCTA=0    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=0    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    **GCTAGC=1**

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=0

# Find the most Common Pattern

AGTCGTCGTTCTCTAAATA**GGCTAGC**

AAATAG=0    AATAGG=0    ATAGGC=0

AGTCGT=0    AGGCTA=0    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=0    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

**GGCTAG=1**



# Find the most Common Pattern

AGTCGTCGTTCTCTAAAT**AGGCTAGC**

AAATAG=0    AATAGG=0    ATAGGC=0

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=0    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

**GGCTAG=1**

# Find the most Common Pattern

AGTCGTCGTTCTCTAAAT**TAGGCT**AGC

AAATAG=0    AATAGG=0    ATAGGC=0

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

**TAGGCT=1**    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

**GGCTAG=1**

# Find the most Common Pattern

AGTCGTCGTTCTCTAA**ATAGGC**TAGC

AAATAG=0    AATAGG=0    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1

# Find the most Common Pattern

AGTCGTCGTTCTCTAA**AATAGG**GCTAGC

AAATAG=0    AATAGG=**1**    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1

# Find the most Common Pattern

AGTCGTCGTTCTCT**AAATAG**GGCTAGC

AAATAG=1    AATAGG=1    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=0

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1

# Find the most Common Pattern

AGTCGTCGTTCTCTAAATAGGCTAGC

AAATAG=1    AATAGG=1    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=0

CGTCGT=0    CGTTCT=0    TAAATA=1

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1

# Find the most Common Pattern

AGTCGTCGTTCTCTAAATAGGGCTAGC

AAATAG=1    AATAGG=1    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=1

CGTCGT=0    CGTTCT=0    TAAATA=1

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=0    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1

# Find the most Common Pattern

AGTCGTCGTTCTCTAAATAGGGCTAGC

AAATAG=1    AATAGG=1    ATAGGC=1

AGTCGT=0    AGGCTA=1    CTAAAT=1

CGTCGT=0    CGTTCT=0    TAAATA=1

TAGGCT=1    TCTAAA=0    TCGTCG=0

TCGTTC=0    TTCTAA=1    GCTAGC=1

GTCGTT=0    GTCGTC=0    GTTCTA=0

GGCTAG=1



Find the most Common Pattern

AGTCGTCGTTCTCTAAATAGGCT  
AGC

Once the first sequence's NUE region is searched, move on to the next sequence.

Searching was so fast the region was expanded to 1->50

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATAT**ATGCTA**

ATGCTA=**1**

TATGCT=0

ATATGC=0

TATATG=0

ATATAT=0

TATATA=0

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATAT**ATGCTA**

ATGCTA=1

TATGCT=1

ATATGC=0

TATATG=0

ATATAT=0

TATATA=0

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATAT**ATATGCTA**

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=0

ATATAT=0

TATATA=0

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=0

TATATA=0

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATAT**ATAT**ATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=**1**

TATATA=0

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=1

TATATA=1

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=2

TATATA=1

CTATAT=0

TCTATA=0

GTCTAT=0



# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=2

TATATA=2

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCT**ATATAT**ATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=**3**

TATATA=2

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=3

TATATA=3

CTATAT=0

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=3

TATATA=3

CTATAT=1

TCTATA=0

GTCTAT=0

# Troubles with Multiple Matches

TGCAGTCAGTCTATATATATGCTA

ATGCTA=1

TATGCT=1

ATATGC=1

TATATG=1

ATATAT=3

TATATA=3

CTATAT=1

TCTATA=1

GTCTAT=0

# Troubles with Multiple Matches

**Biased towards repeating pattern**

**Added code to count patterns only once per sequence.**

# Results – 6 character length

## Multi-count

TTTTTT	3394
AATAAA	2242
TTTGTT	2224
TTTTGT	2122
TTGTTT	2025
AAAAAA	2006
ATATAT	1945
TGTTTT	1937
ATTTTT	1897
TTTCTT	1858

## Single-count

AATAAA	2021
TTTGTT	2021
TTTTGT	1993
TTGTTT	1859
TGTTTT	1806
ATTTTT	1795
TTTCTT	1720
TTTTAT	1695
TATTTT	1662
TTTTTT	1629

# Results – 7 character length

## Multi-count

TTTTTTT	1717
TTTGTTT	1043
AATAAAA	988
AAAAAAA	982
TTTTGTT	980
TATATAT	951
ATATATA	895
AAATAAA	835
TTGTTTT	832
TTTCTTT	803

## Single-count

TTTGTTT	972
AATAAAA	945
TTTTGTT	940
TTTTTTT	819
TTGTTTT	796
TTTTCTT	767
TTTCTTT	758
AAATAAA	757
TTTTTGT	724
TTTATTT	716



# Results – 8 character length

## Multi-count

TTTTTTTT	898
ATATATAT	492
TATATATA	479
AAAAAAAA	475
TTTTGTTT	459
TTTGTTTT	447
AAATAAAA	363
TTTTTCTT	361
TTTTCTTT	349
TTTTTGTT	338

## Single-count

TTTTGTTT	436
TTTGTTTT	429
TTTTTTTT	406
TTTTTCTT	357
AAATAAAA	347
TTTTCTTT	341
TTTTTGTT	337
ATTTTTTT	326
TTTCTTTT	322
TAATAAAA	319

# Results – 9 character length

## Multi-count

TTTTTTTTTT	499
TATATATAT	290
ATATATATA	274
AAAAAAAAAA	229
TTTTGTTTT	206
TTTGTTTTT	169
TTTTTCTTT	164
ATAAATAAA	158
TTTTCTTTT	158
TTTTTGTTT	155

## Single-count

TTTTTTTTTT	234
TTTTGTTTT	194
TATATATAT	177
TTTGTTTTT	168
TTTTTCTTT	162
ATATATATA	161
TTTTCTTTT	155
ATTTTTTTT	154
TTTTTGTTT	154
ATAAATAAA	146

# Results – 10 character length

## Multi-count

TTTTTTTTTT	273
ATATATATAT	178
TATATATATA	171
AAAAAAAAAA	90
AATAAATAAA	83
ATTTTTTTTT	82
TTTTTTTTTG	81
TTTGTTTTT	79
TTTTTTTCTT	78
TTTTCTTTT	74

## Single-count

TTTTTTTTTT	123
ATATATATAT	93
TATATATATA	90
ATTTTTTTTT	82
TTTTTTTTTG	81
TTTGTTTTT	79
TTTTTTTCTT	78
AATAAATAAA	76
GTTTTTTTTT	74
TTTTCTTTT	73

# Results – 2D Matrix

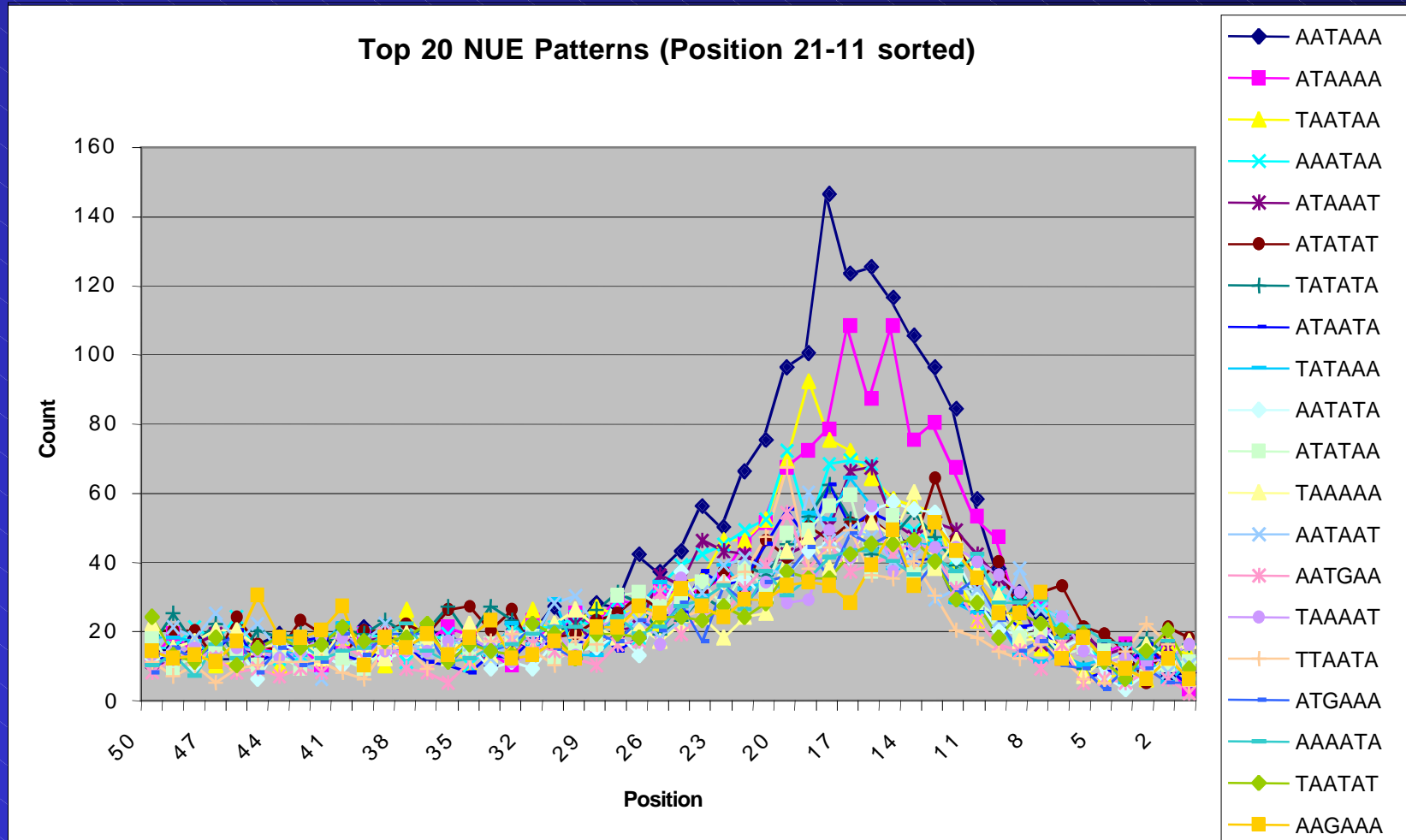
**Needed more information about the locations**

**Expanded the pat\_freq array into a 2D matrix**

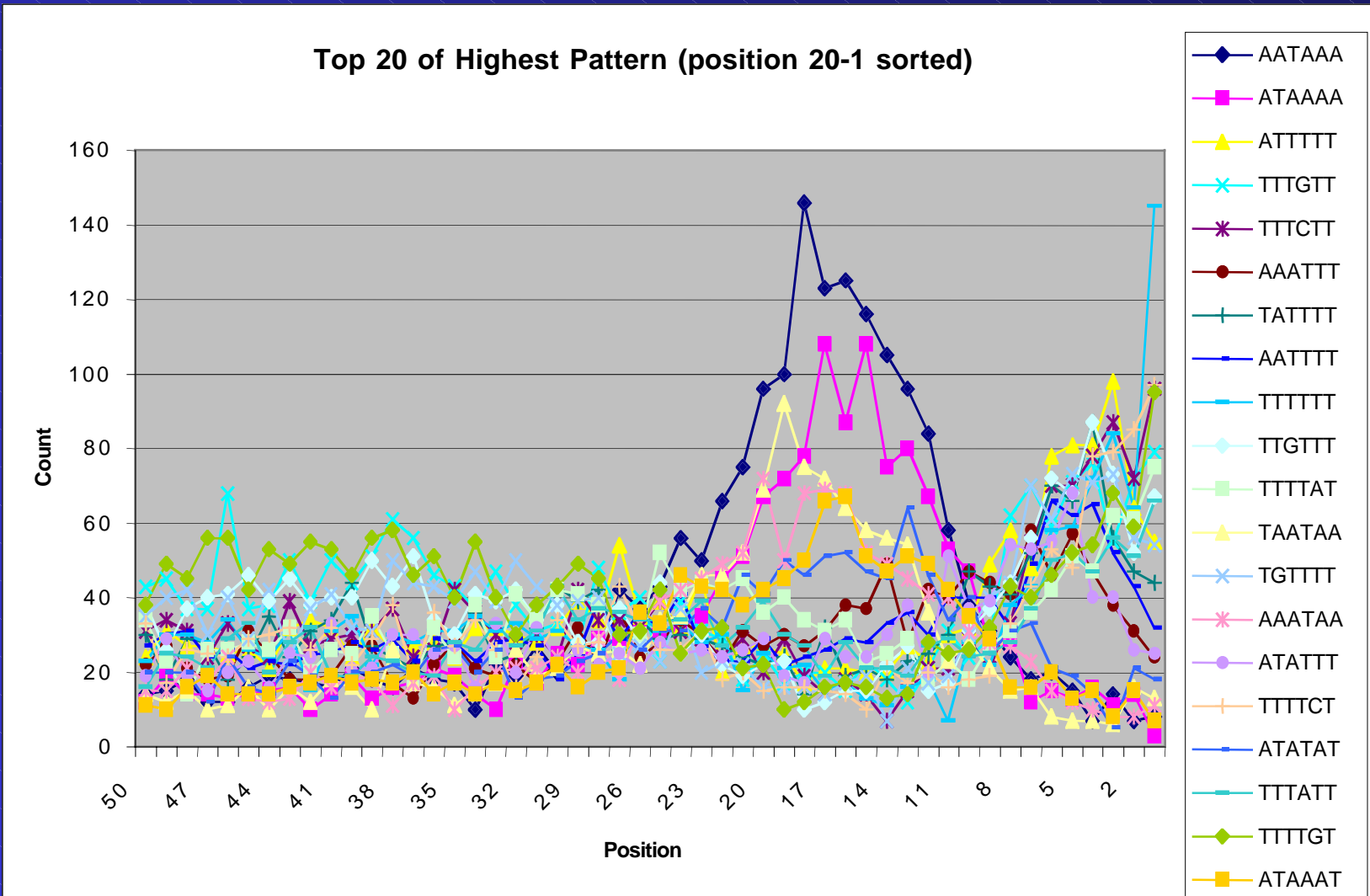
**Now we can count the most common pattern  
on a location basis**

**Graphs show frequency by location**

# Summary Findings



# Summary Findings



# Conclusions

- Search revealed the NUE patterns that are in the expected range of location
- Discovered unknown consensus of patterns near the cleavage site

# Next step

- Expand the search to the FUE region
- Tie the NUE and FUE searches together
- Package the program so it more general purpose for others to use
- Install at OSC



# Acknowledgements

- Dr. Brain Haas (TIGR) for providing raw databases
- Cray Inc, for SV1ex computer time
- OSC for SV1 computer time and support