

The Role of Bioinformatics in Gene Discovery

James J. Vincent

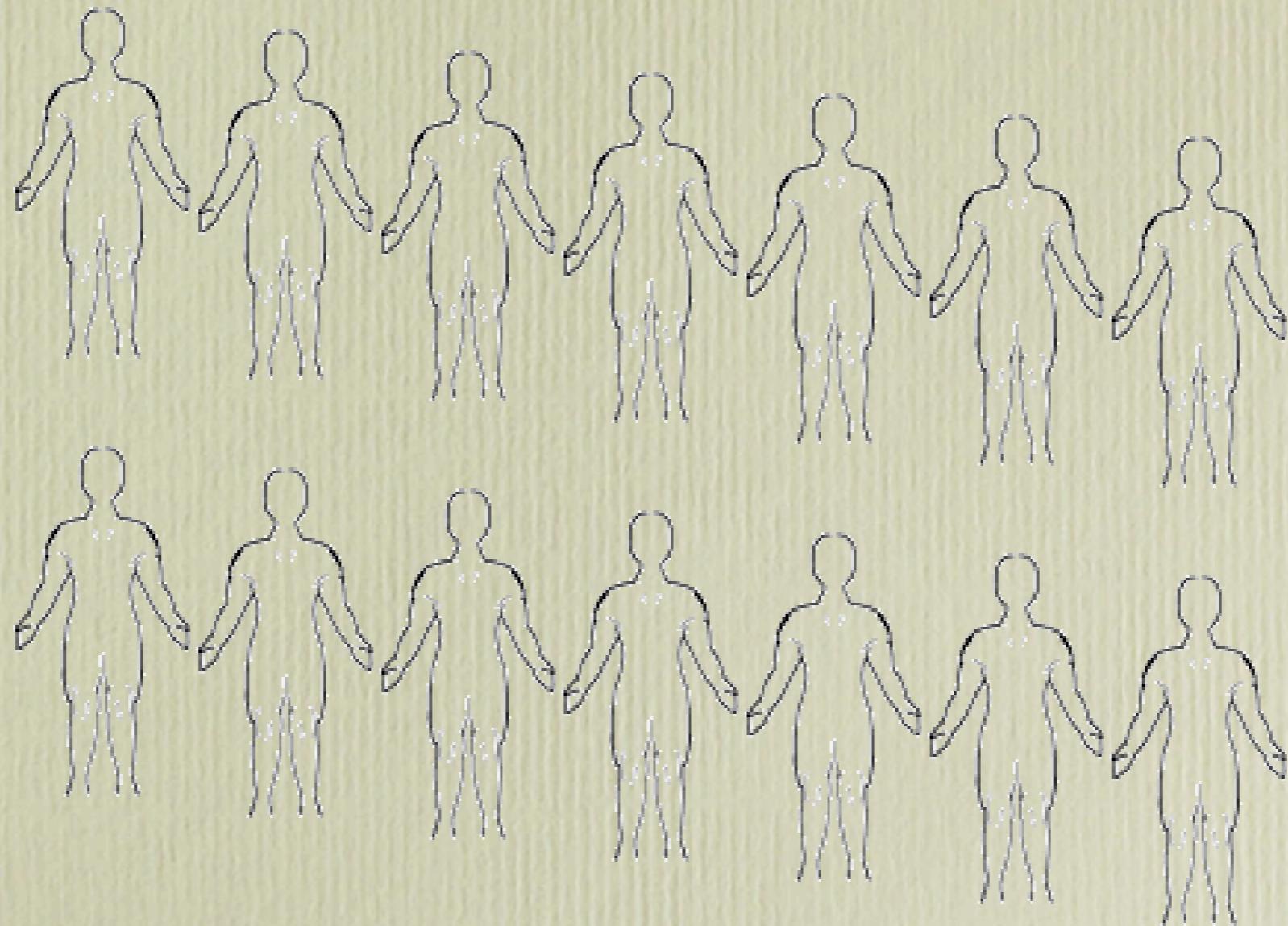
Laboratory of Molecular Biology
National Cancer Institute
National Institutes of Health



Overview

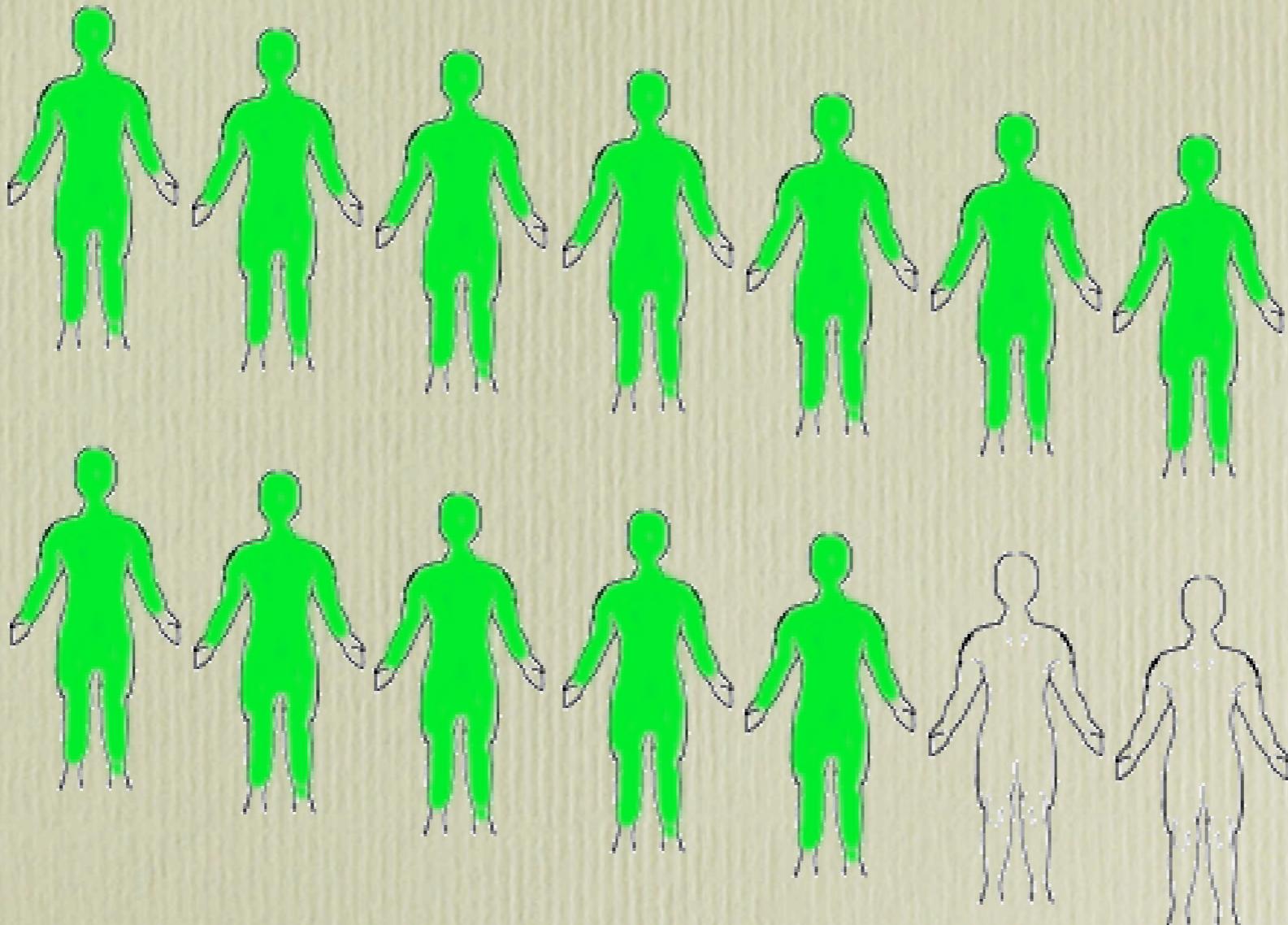
- Why we look for new genes
- How we find them with a computer
- We found what we were looking for

We make immunotoxins ...



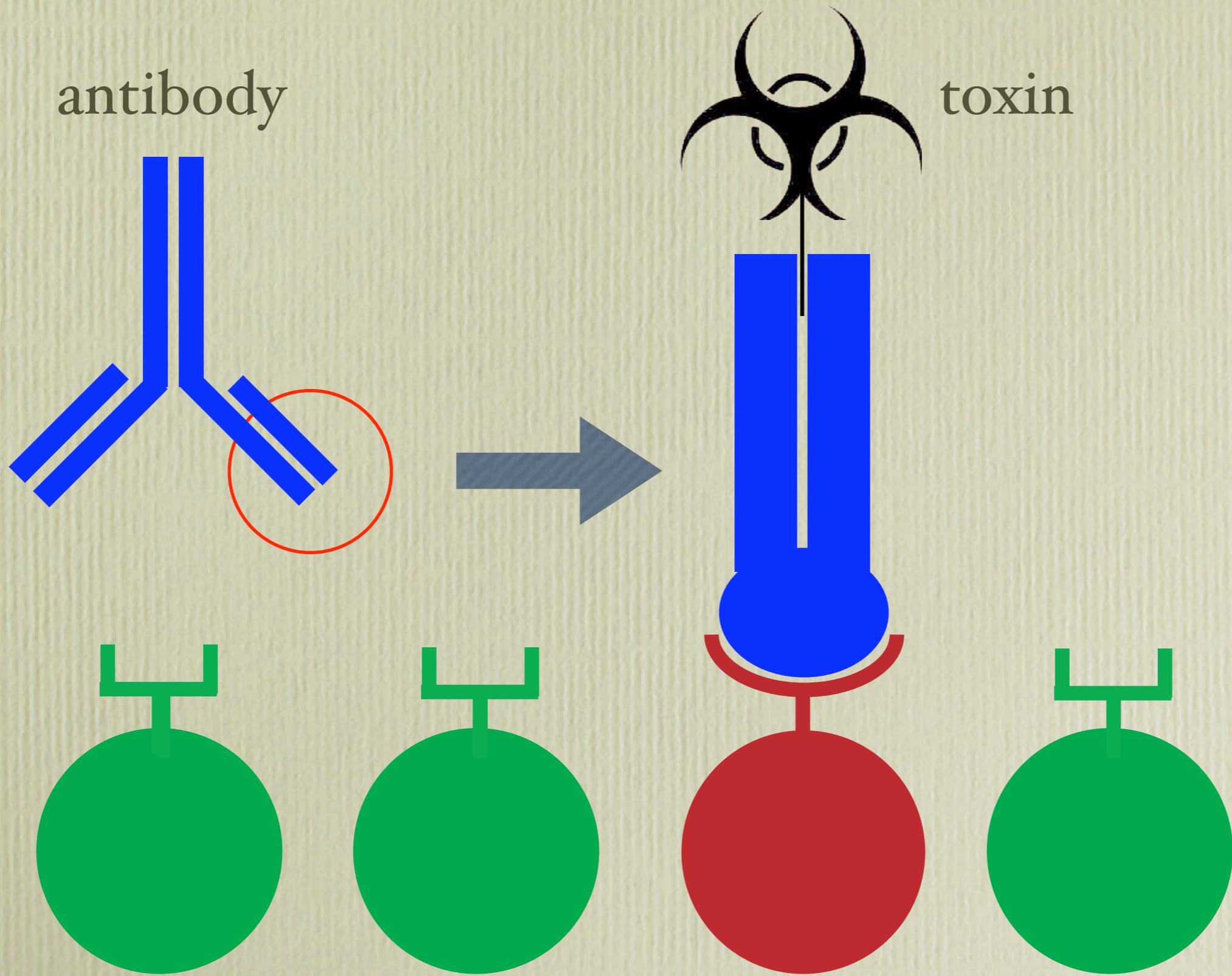
14 Patients with Hairy Cell Leukemia
Phase I Clinical Trial

... and they work well.

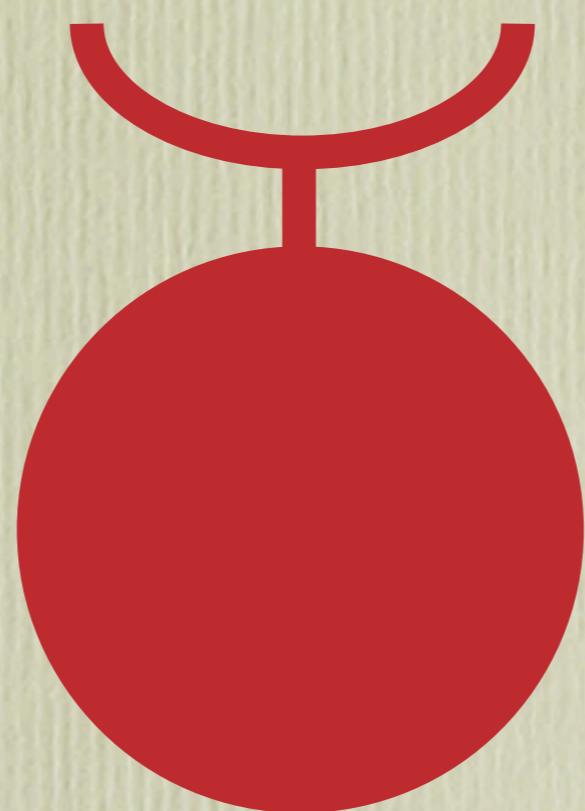


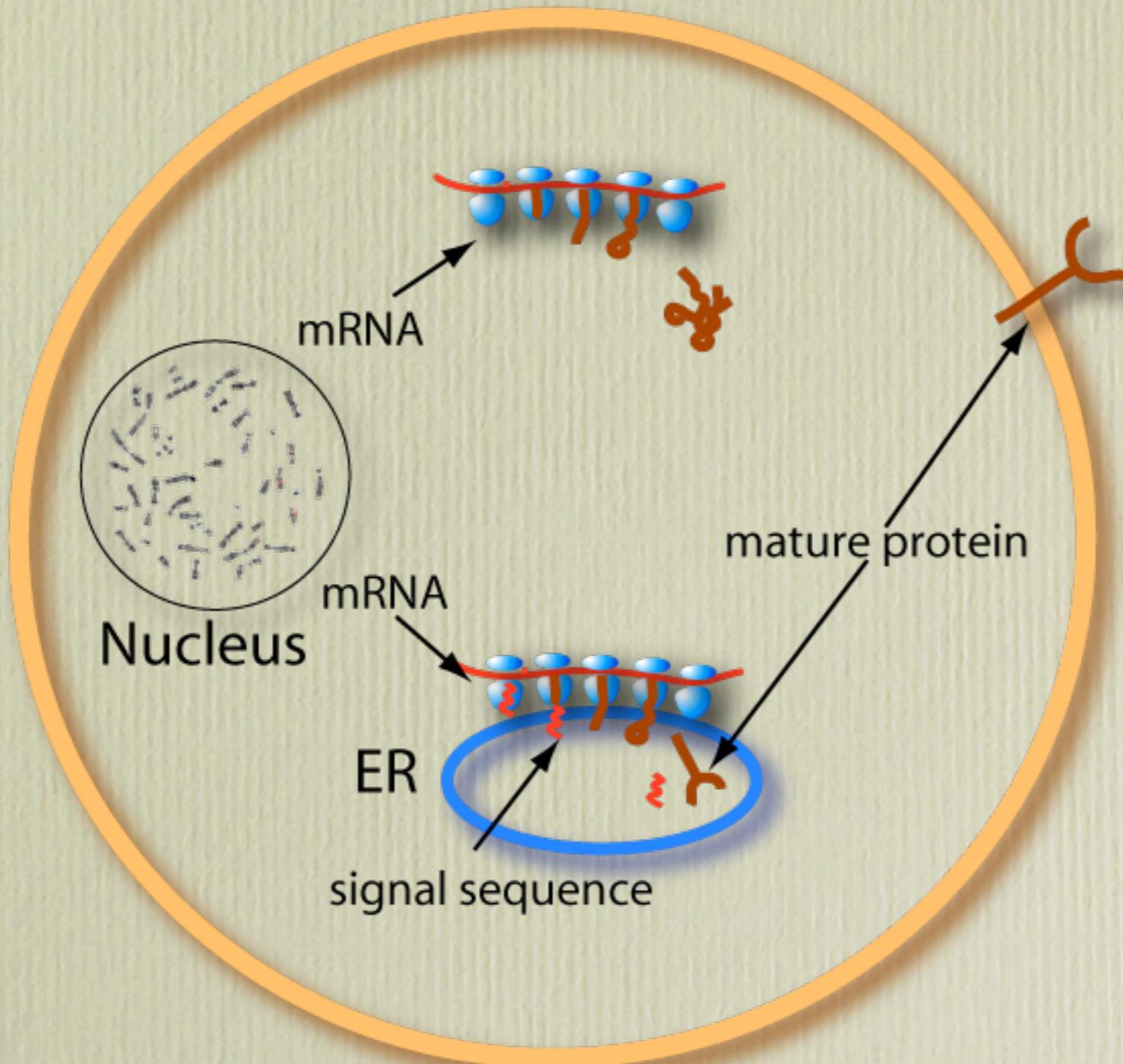
86% Complete Remission

How We Make Immunotoxins



Find the Receptor



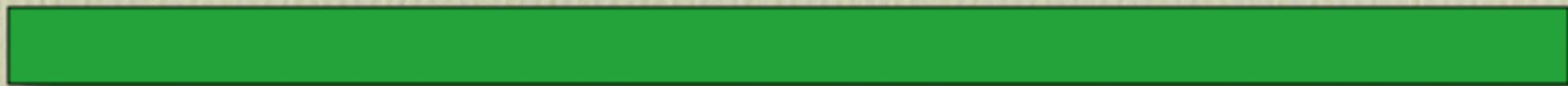


The Human Genome

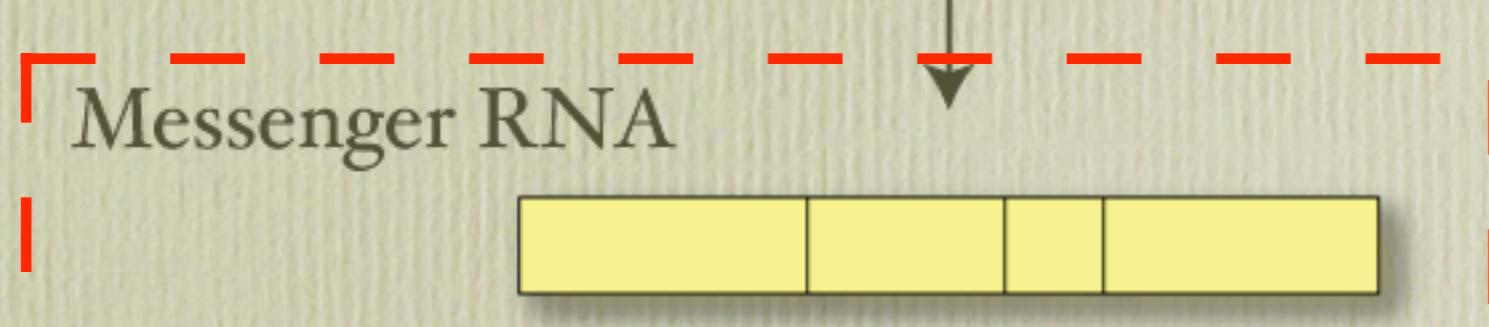
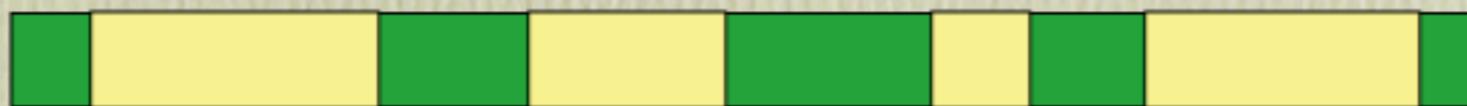
Mostly Finished

The screenshot shows the NCBI genome browser interface for *Homo sapiens* build 33. The top navigation bar includes links for PubMed, Nucleotide, Protein, Genome, Structure, and PopSet, along with search fields for "Search for", "on chromosome(s)", and "haplotype", and a "Find" button. Below the search bar are buttons for "Show linked entries", "Help", "FTP", and "MVhome". The main content area displays a chromosome map with chromosomes 1 through 22, X, Y, and MT labeled. Above the chromosomes, a genomic track shows DNA sequence with various genes and regulatory elements highlighted in different colors. The genes *gjfc*, *thrL*, *cynR*, and *ushfA* are specifically labeled.

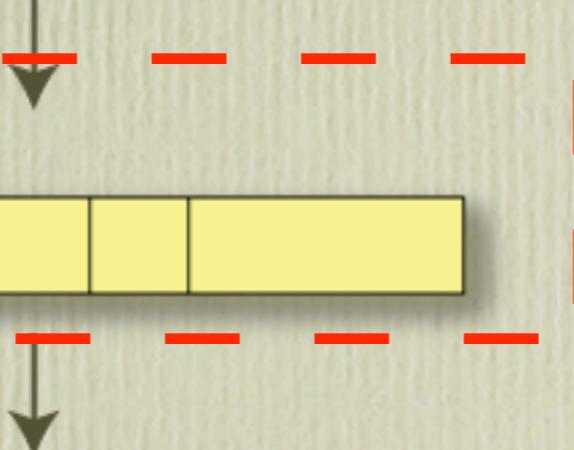
Chromosomal DNA



Nuclear RNA



Messenger RNA

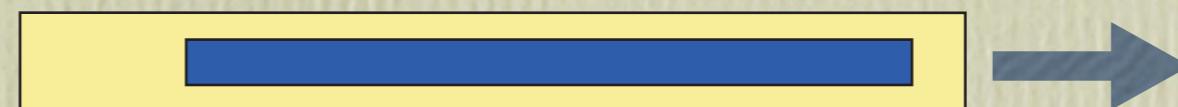


Protein

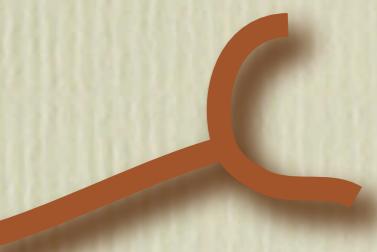
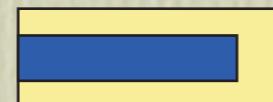


Expressed Sequence Tags

Messenger RNA



Expressed Sequence Tag



Protein

EST

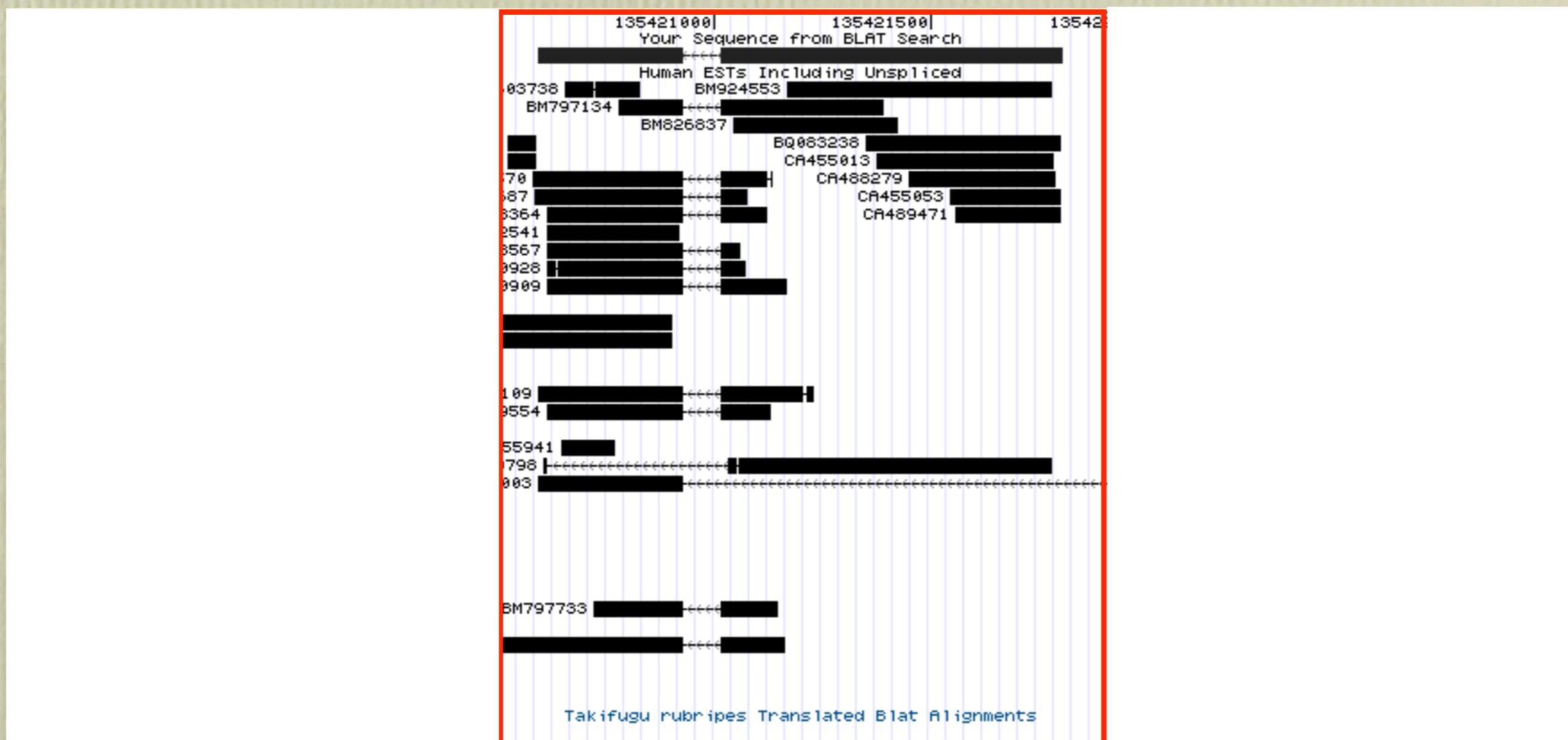


dbEST

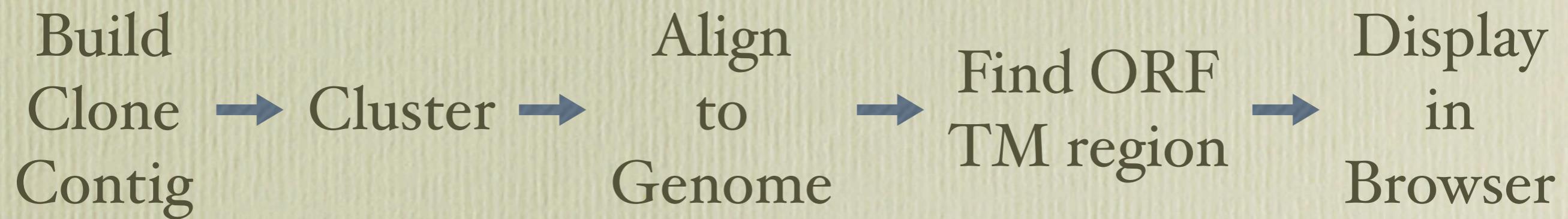
- 5 million ESTs
- 8000 libraries
- 1500 tissue types

Cluster together into genes

Clustering Results



Sequence Analysis Pipeline



The Value of a C.S. Degree



Script web tools



Large relational databases



Reusable software objects



Software for other people



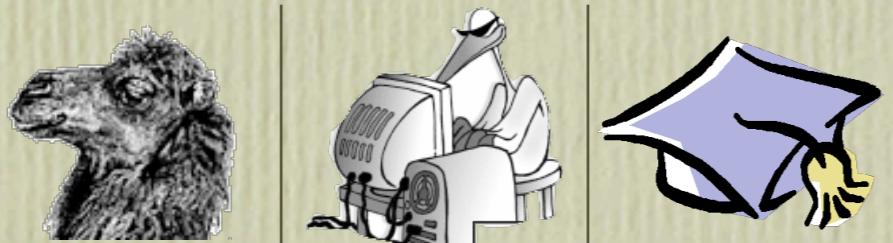
Software for other software



Robust, documented software



System Administration Backup Networking Queueing System



Breast Cancer Genes from dbEST

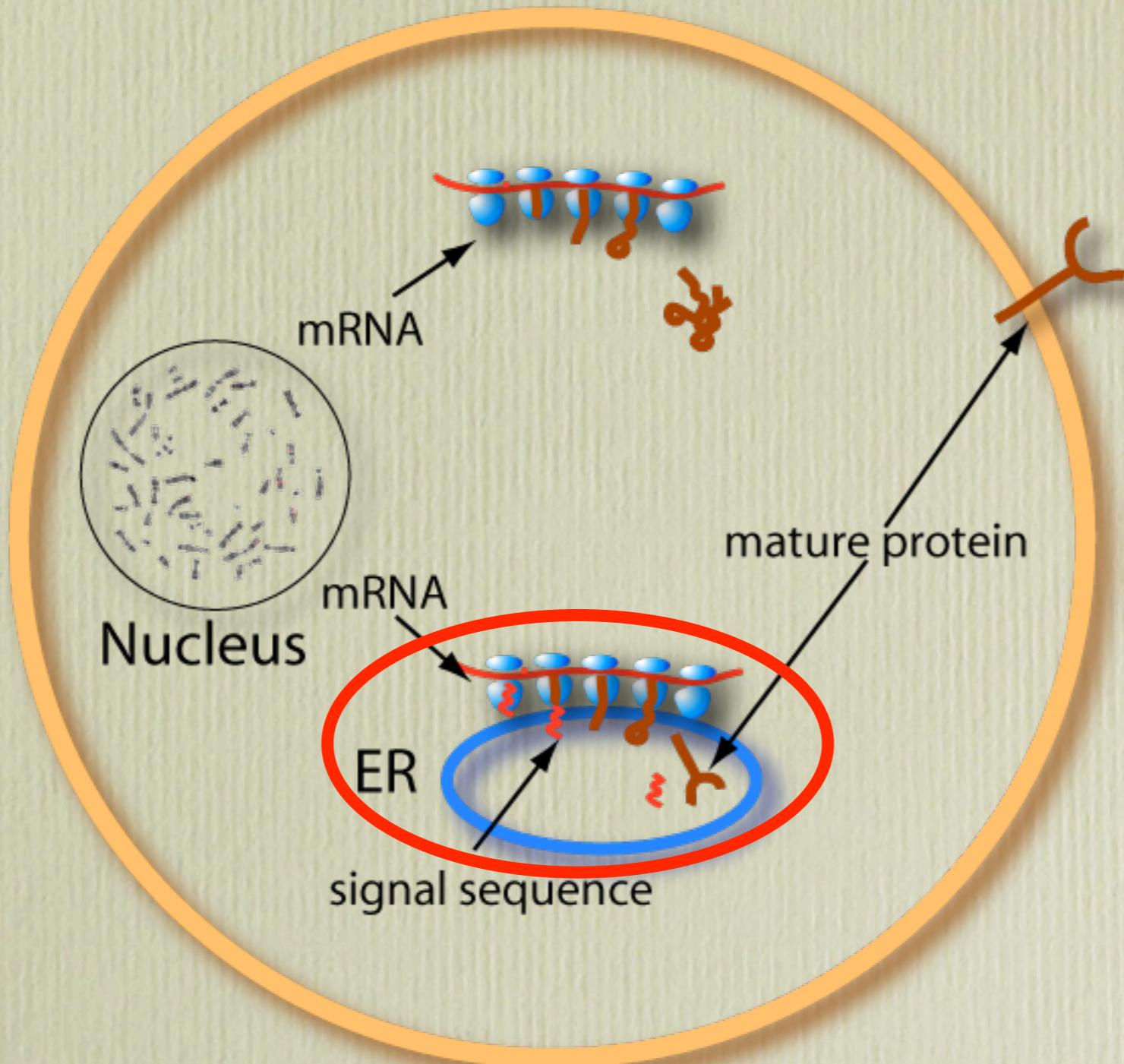
O

Kristi Egland, PhD

Molecular Biologist

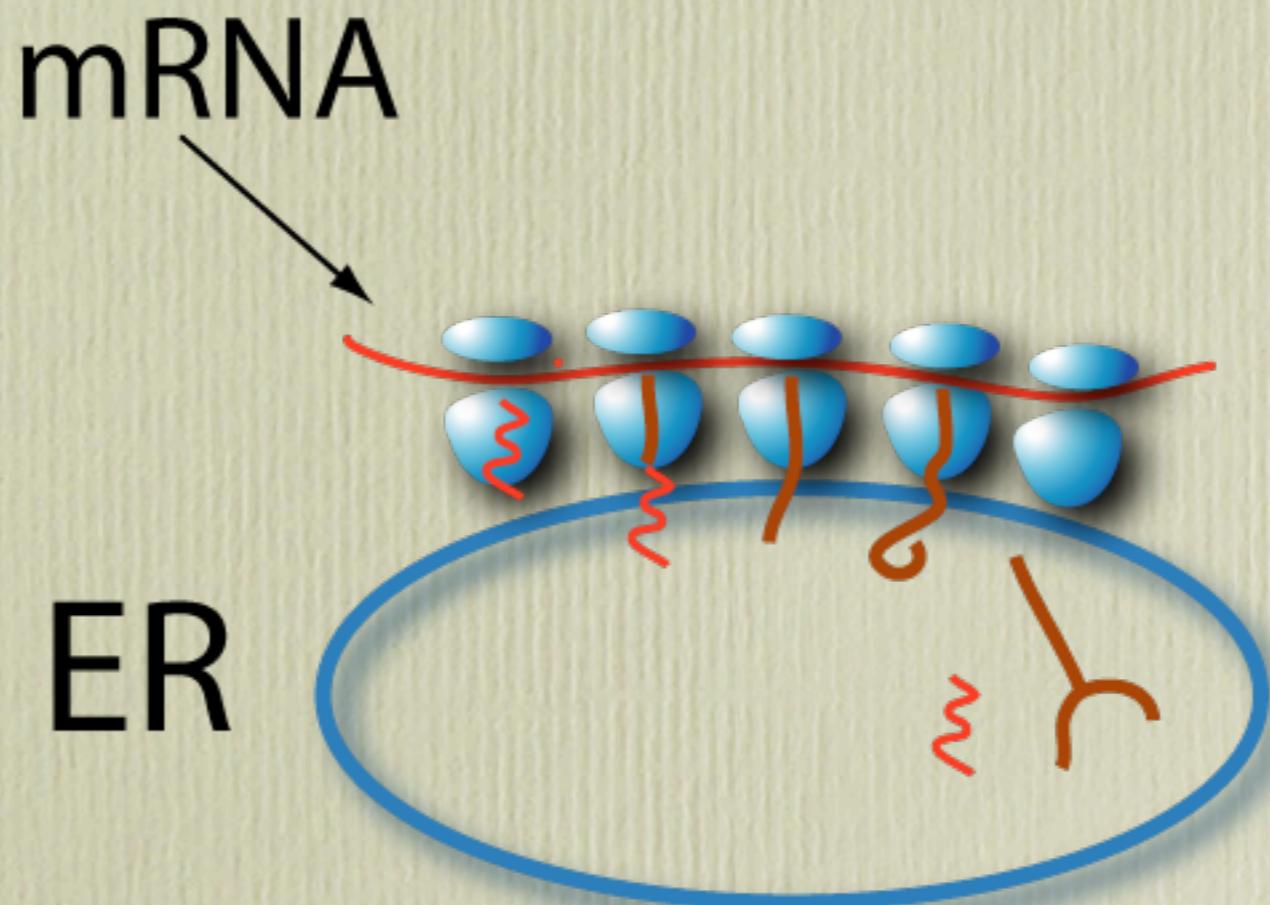


We Make Our Own ESTs



The MAPcL

Membrane Associated Polyribosomal cDNA Library

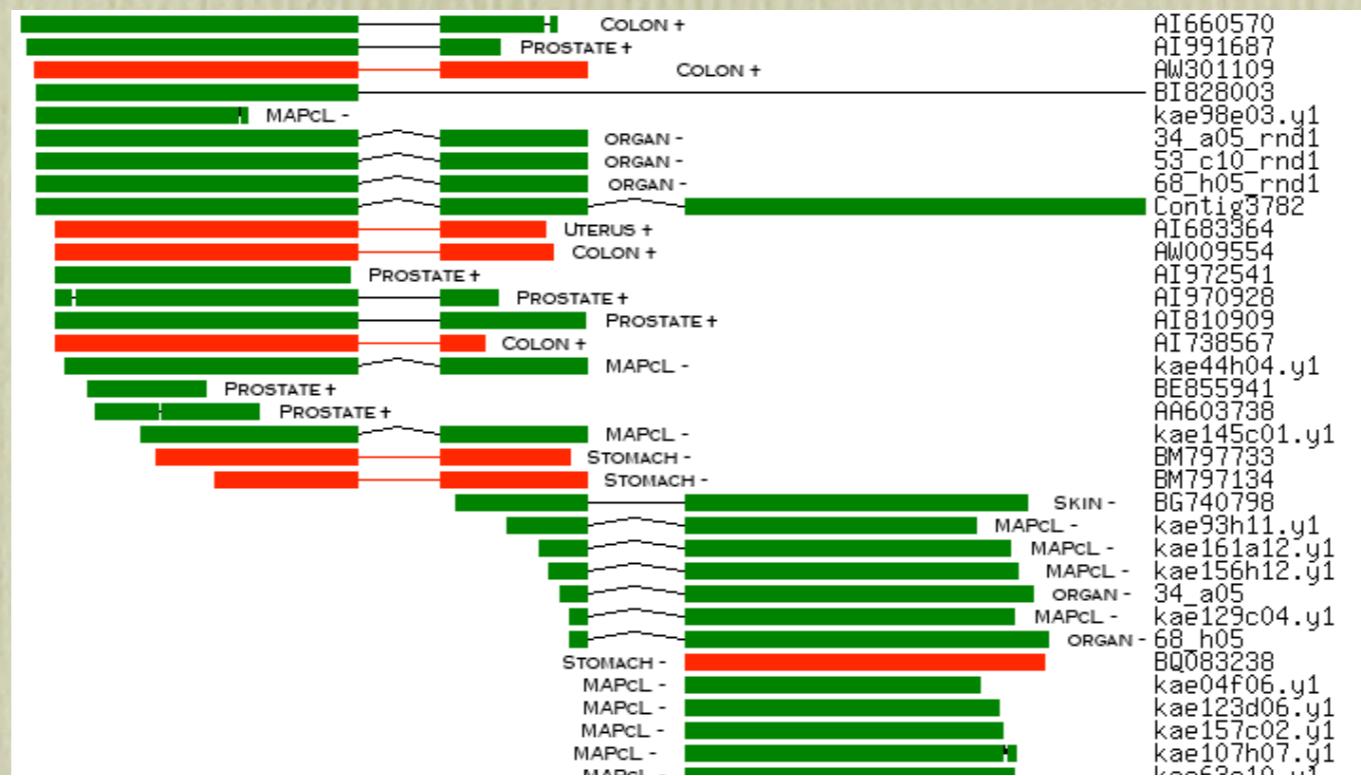


Two years + \$\$\$\$\$ → 25,000 sequences

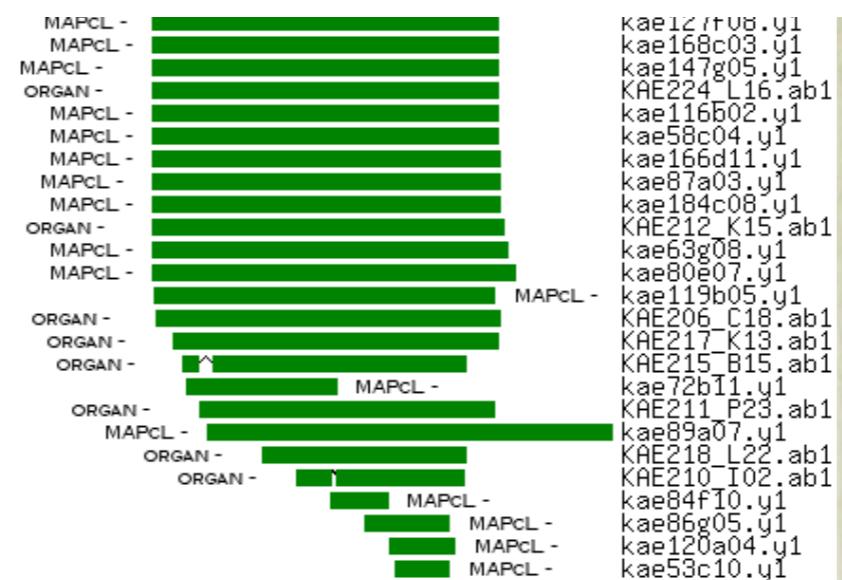
The MAPcL is Above Average

- Enriched for membrane protein genes
- Genes from critical organs removed
- Sequenced from the 5' end

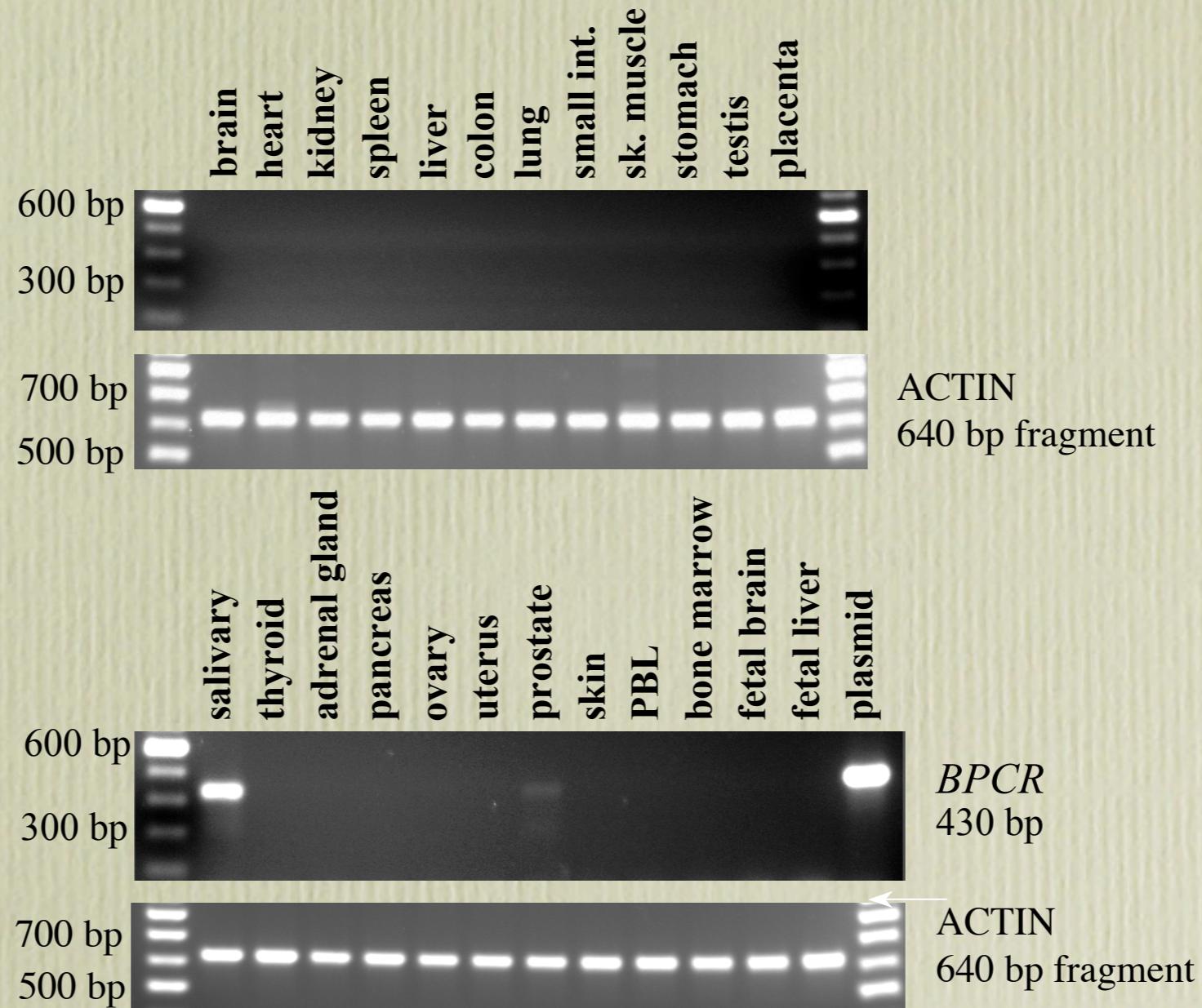
Breast Cancer Genes from MAPcL



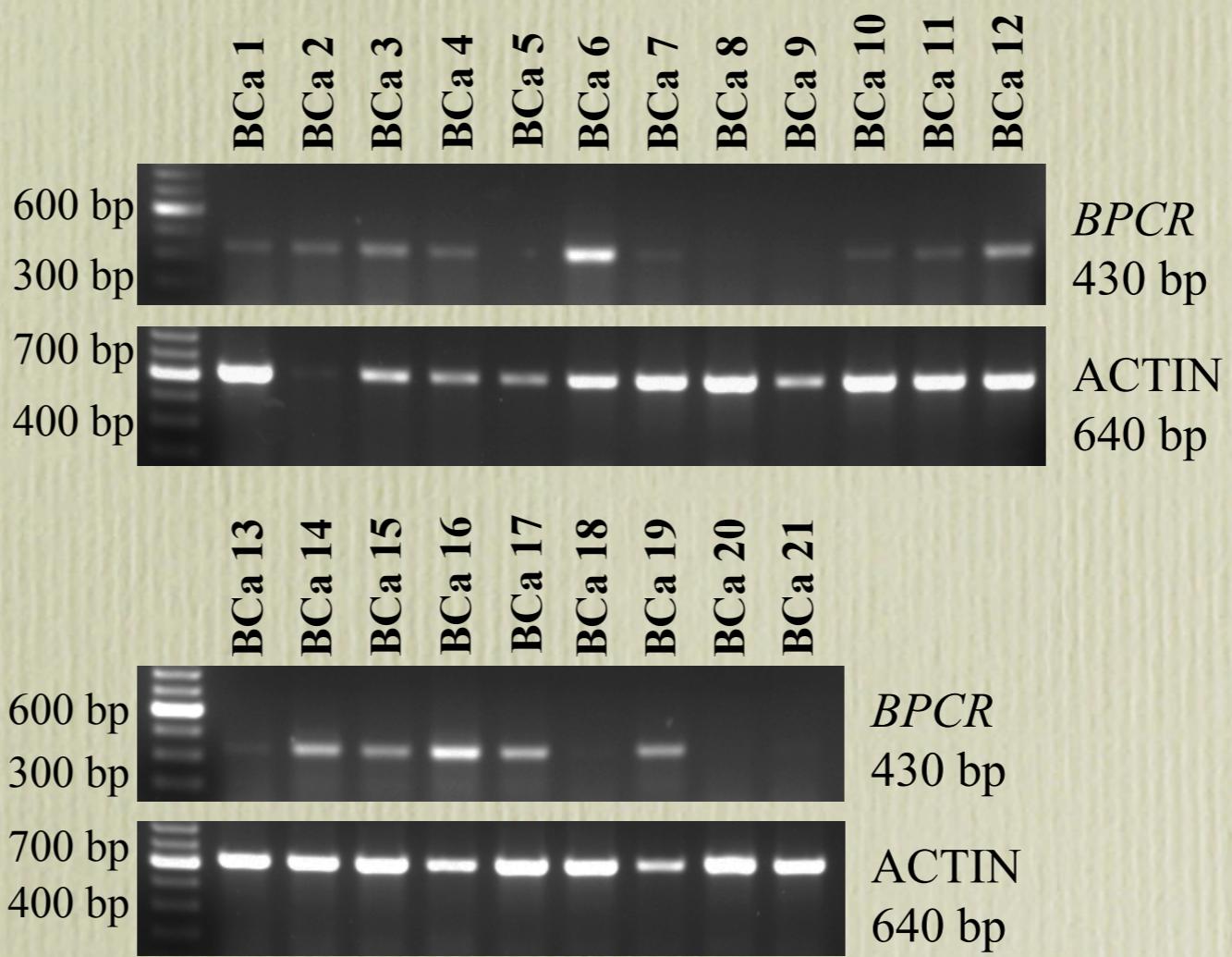
Breast Prostate and Salivary Gland Receptor - BPSR



RT-PCR Analysis of *BPCR* in Normal Tissues

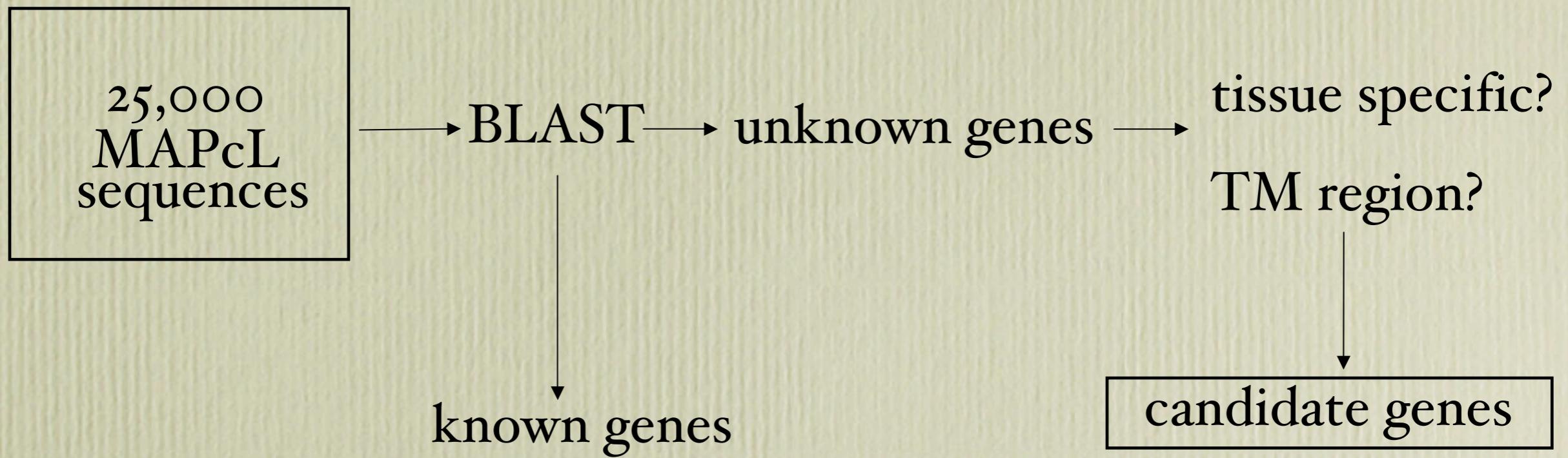


RT-PCR Analysis of *BPCR* Expression in Breast Cancers



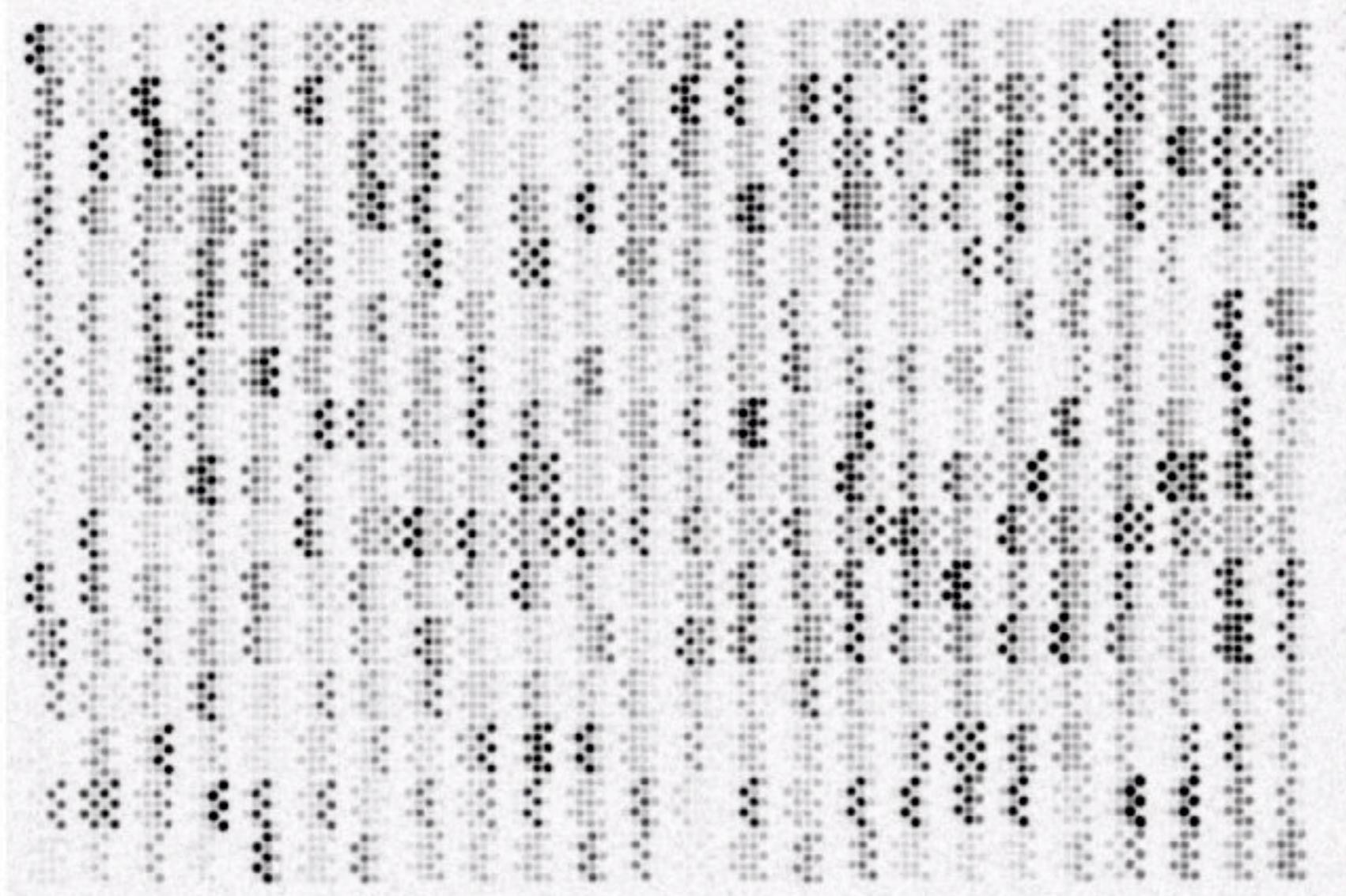
Analysis of MAPcL Sequences

How do we find new targets?



1400 unknown genes selected for array analysis

All at once...
1400 genes on one array



Sample 1

Summary

- We found what we were looking for
- High quality software plays a critical role
- C.S. degree directly relevant to life sciences research

Kristi Egland
Ira Pastan
Bob Strausberg
Kevin Becker





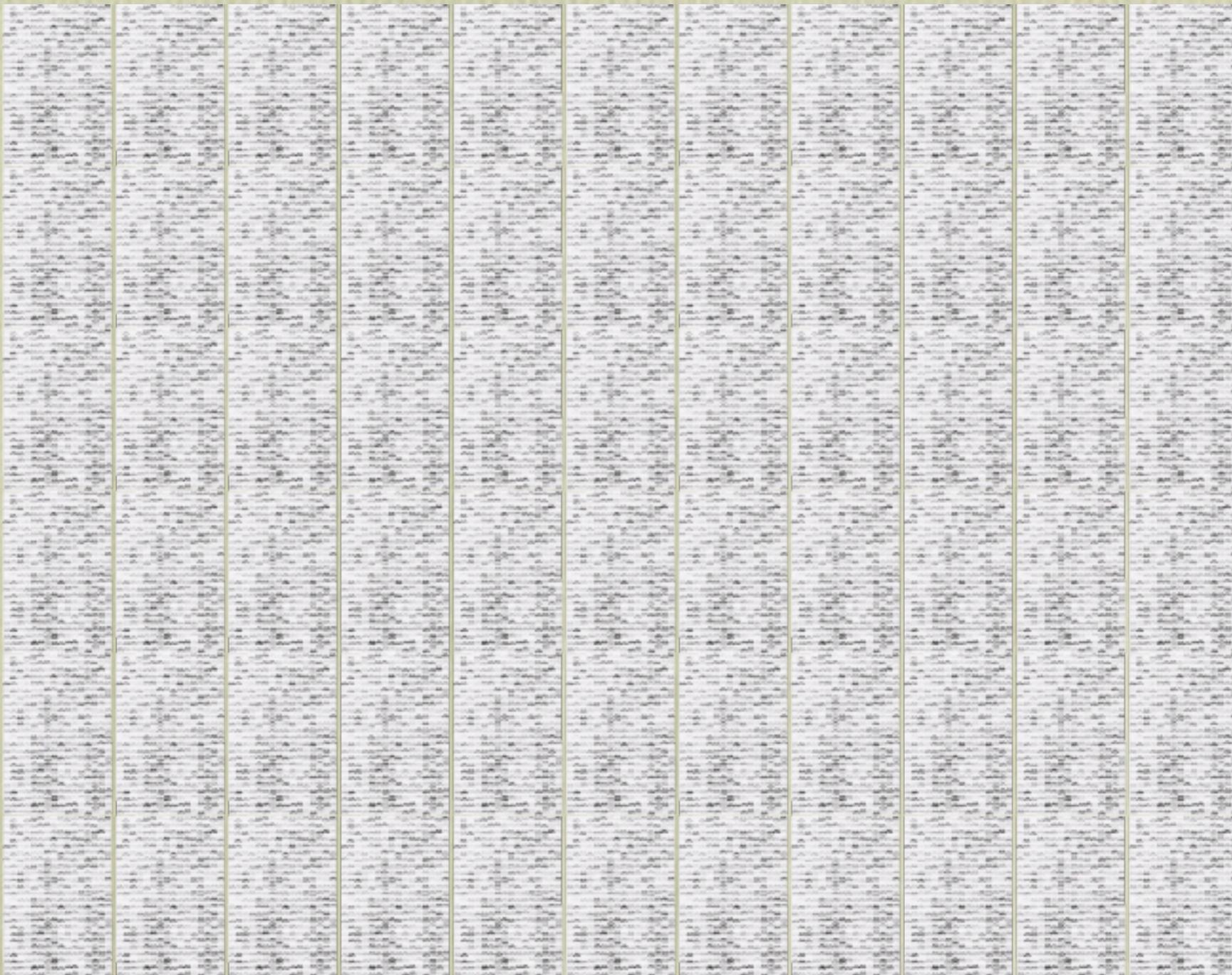


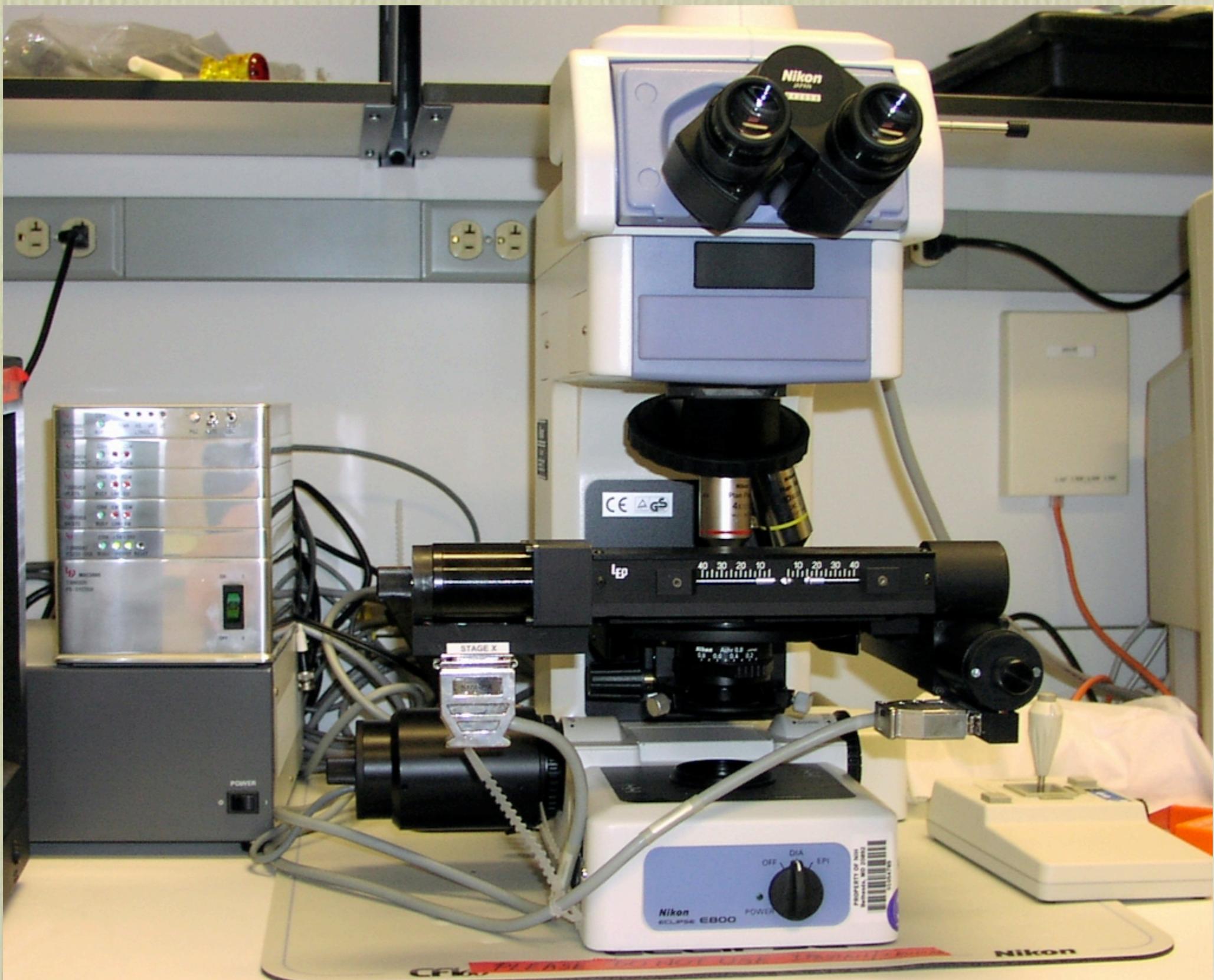


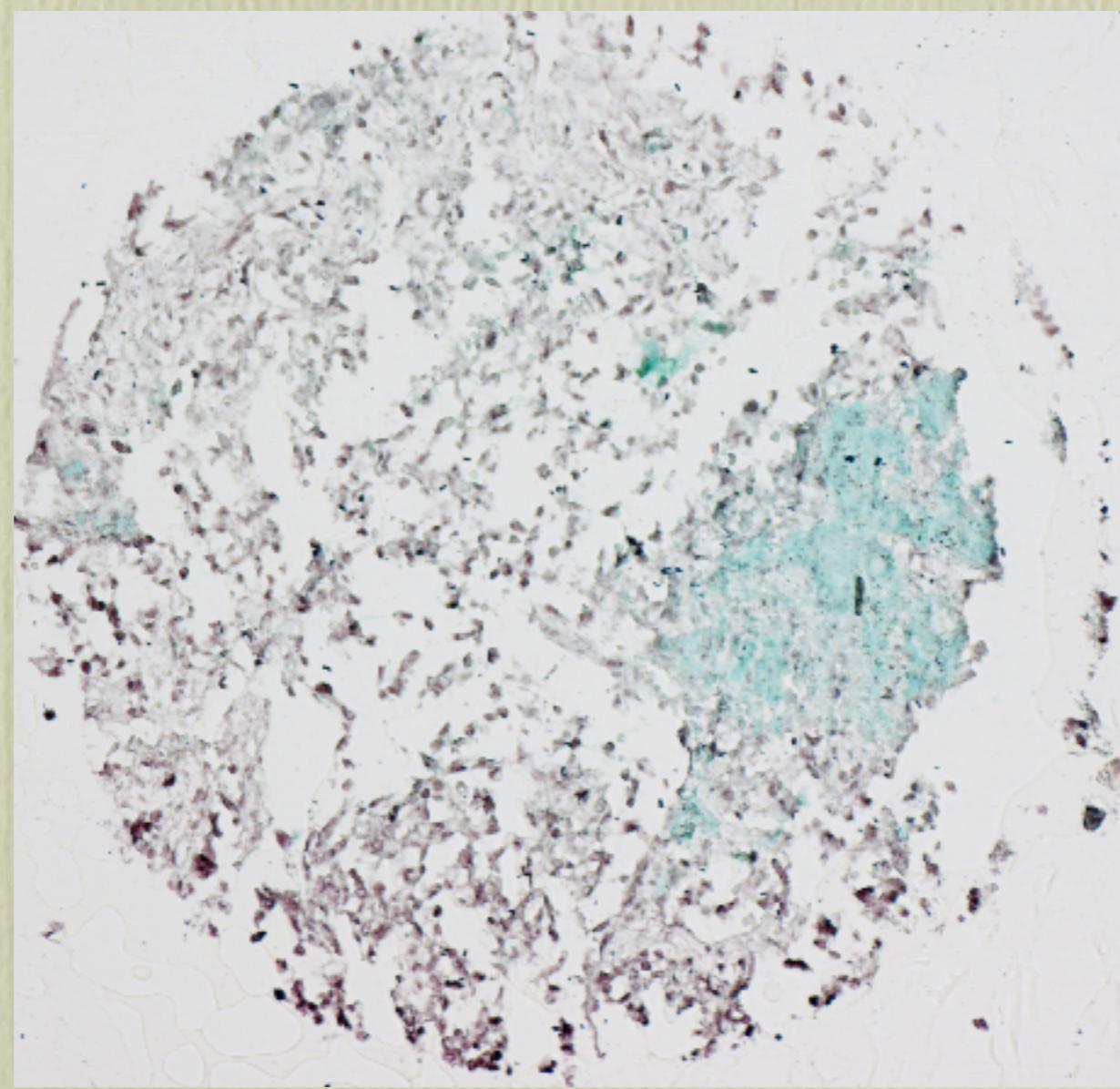


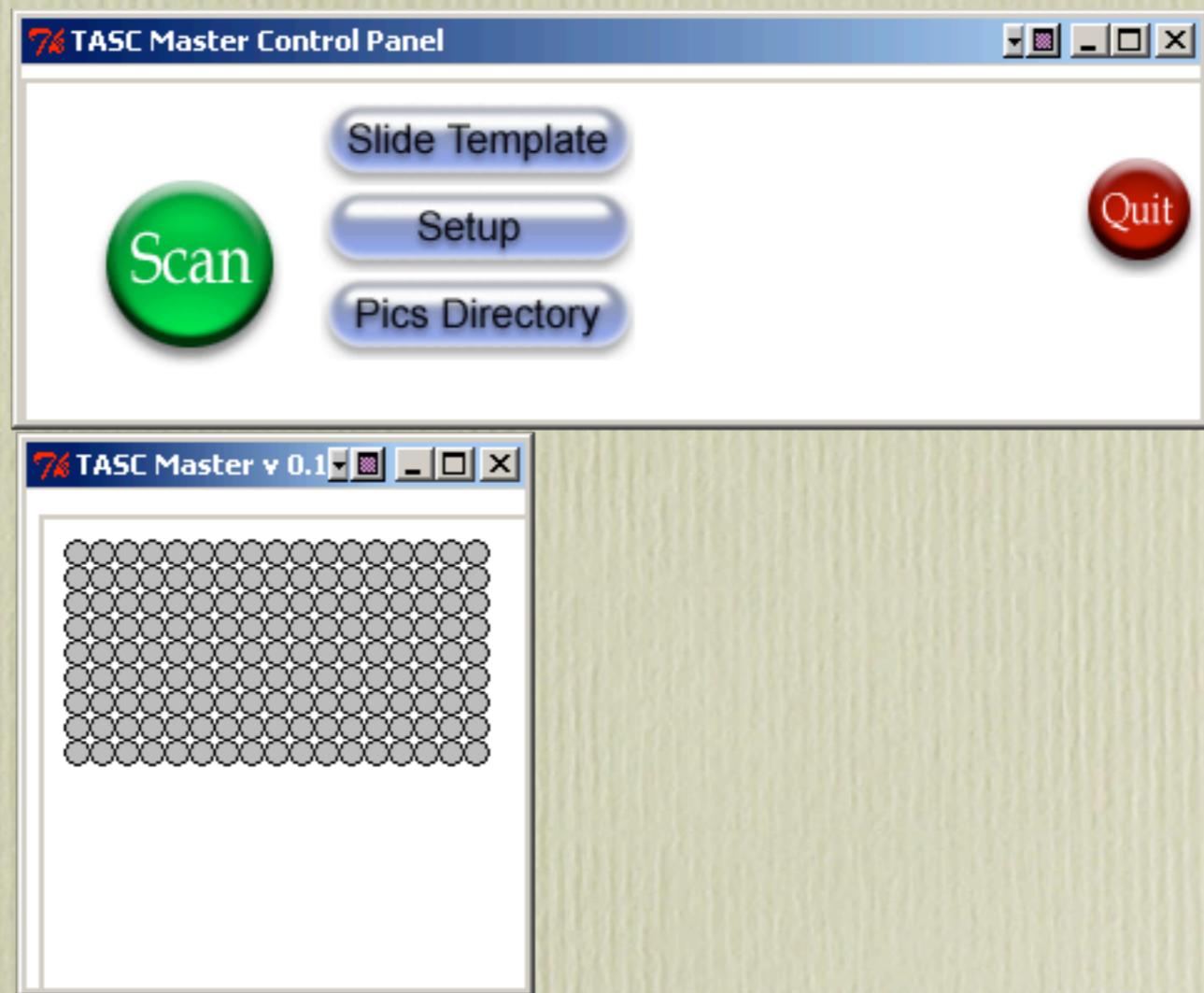
Too Much

134 arrays x 10,000 spots = 1.3M data points









- RefSeq array - expression to pathway
- dbEST zScore compare to RefSeq array
- Expression browser - all info
- Array experiment interface for public databases
- Lit search tool hooked to othe rdatasbases
- Promoter sequence extract, analyze, compare