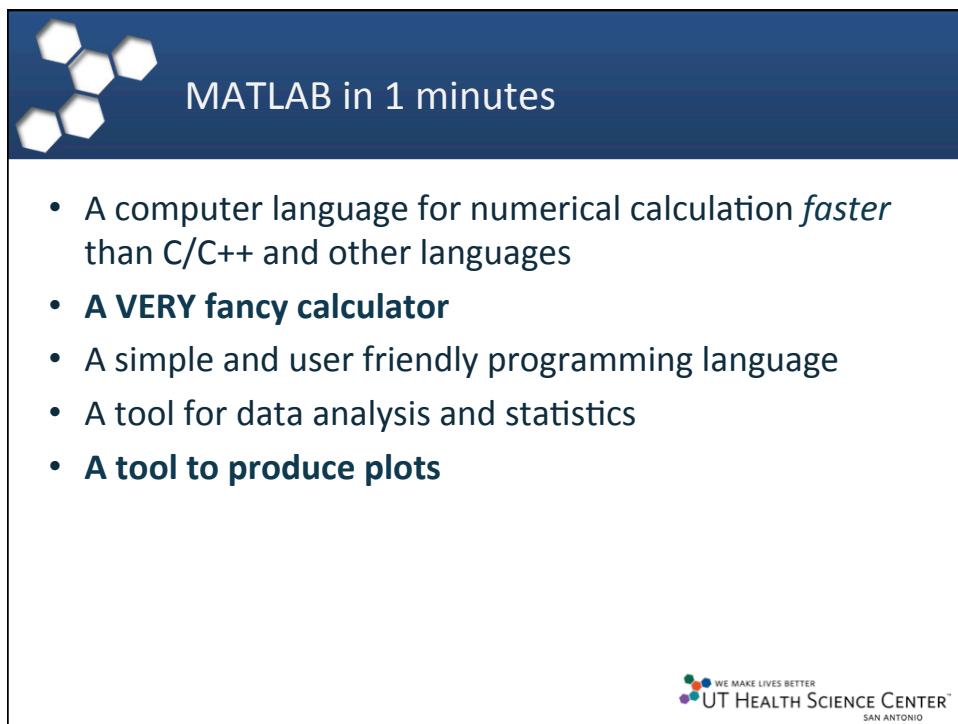


A presentation slide with a dark blue header and a white footer. The header contains the title 'MATLAB & Microarray Data Analysis' in white, along with the name 'Yidong Chen', 'Computational Biology and Bioinformatics', and an email address 'cheny8@uthscsa.edu'. The footer features the UT Health Science Center logo with the tagline 'WE MAKE LIVES BETTER' and 'SAN ANTONIO'.

## MATLAB & Microarray Data Analysis

Yidong Chen  
Computational Biology and Bioinformatics  
cheny8@uthscsa.edu

WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER™  
SAN ANTONIO



A presentation slide with a dark blue header and a white footer. The header contains the title 'MATLAB in 1 minutes' and the UT Health Science Center logo. The main content area lists several bullet points about MATLAB's capabilities.

## MATLAB in 1 minutes

- A computer language for numerical calculation *faster* than C/C++ and other languages
- **A VERY fancy calculator**
- A simple and user friendly programming language
- A tool for data analysis and statistics
- **A tool to produce plots**

WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER™  
SAN ANTONIO



## Statistics Toolbox

**Apply statistical algorithms and probability models**

**Key Features**

- Regression techniques, including linear, generalized linear, nonlinear, robust, regularized, ANOVA, and mixed-effects models
- Univariate and multivariate probability distributions
- Random and quasi-random number generators and Markov chain samplers
- Hypothesis tests for distributions, dispersion, and location, and design of experiments (DOE) techniques
- Supervised/Unsupervised machine learning algorithms, including support vector machines (SVMs), k-means and hierarchical clustering, Gaussian mixtures, and hidden Markov models

 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



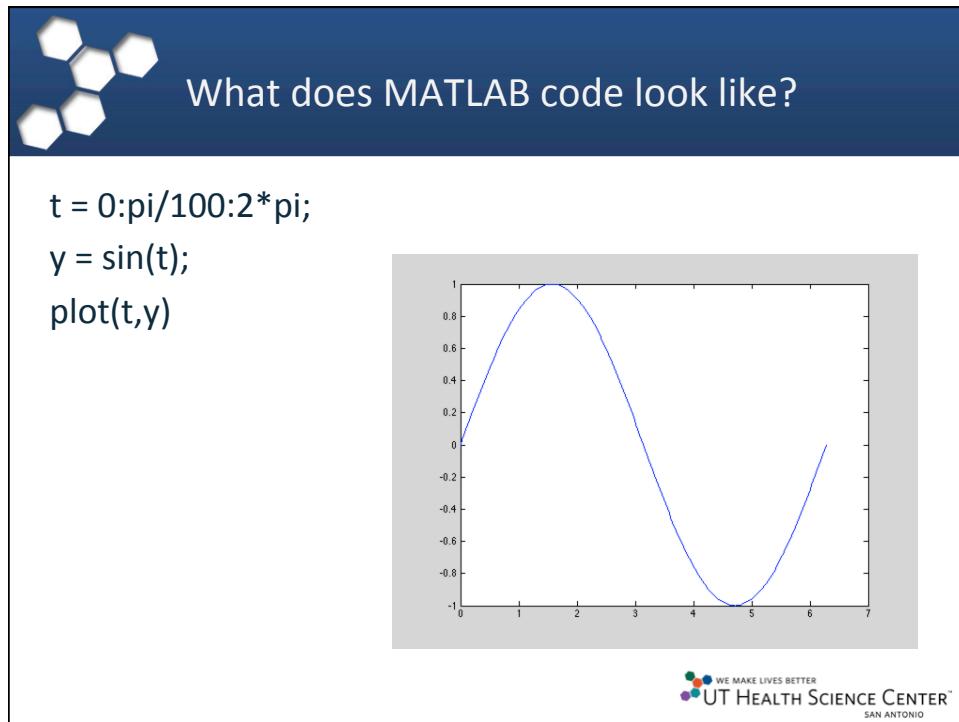
## Bioinformatics Toolbox

**Read, analyze, and visualize genomic, proteomic, and high throughput profiling data**

**Key Features**

- **Next Generation Sequencing** analysis and browser
- Sequence alignment analysis and visualization
- **Microarray**, including reading, filtering, normalizing, and visualization
- **Mass spectrometry**, preprocessing, classification, and marker identification
- **Phylogenetic tree analysis**
- Graph theory functions, including interaction maps, hierarchy plots, and pathways
- Data import from genomic, proteomic, and gene expression files, including SAM, FASTA, CEL, and CDF, and from databases such as NCBI and GenBank

 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



What does MATLAB code look like?

```
t = 0:pi/100:2*pi;
y = sin(t);
plot(t,y)
```

A plot of the sine function from 0 to 2π, showing one full cycle of the wave.

WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



MATLAB Desktop (Mac, R2012a)

File Edit Debug Parallel Desktop Window Help

/Users/cheny8

Shortcuts How to Add What's New

Current Folder

Name

- zhujack\_2008-12-22.sql
- zhujack\_2008-12-22 2.sql
- weka.log
- tmp.txt
- tmp.err
- test 2.fp7
- test.fp7
- Table S1.txt
- SqlViewerHistory.props
- Send Registration 2
- Send Registration
- Research Statement\_Degeng 2...
- Research Statement\_Degeng.doc
- Rao\_IPA\_1a\_obs5.txt
- Rao\_IPA\_1a\_obs4.txt
- Rao\_IPA\_1a\_obs3.txt
- Rao\_IPA\_1a\_obs2.txt
- Rao\_IPA\_1a\_obs1.txt
- Project\_GGT112344.tar.bz2
- passport\_YidongChen.pdf
- NucleotidePoolStudy\_revYidong...
- Net2\_53genes\_CellGrowth\_Jar...
- Net1\_53genes\_NFkB\_largeNet...
- mockup.pptx
- matlab\_crash\_dump.64228-1

Your Directory

Command Window

Enter your command here

Workspace

Your variable

Start Ready

WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



## Some simple functions

```

>> A = [1 3 2 1; 5 10 11 8; 9 6 7 1; 4 15 14 1]
MATLAB displays the matrix you just entered.
A =
    1 3 2 1
    5 10 11 8
    9 6 7 1
    4 15 14 1
>> sum(A)           >> mean( A )
ans =
    19   34   34   11
                                         ans =
                                         4.7500  8.5000  8.5000  2.7500
>> B = sum(A')
B =
    7   34   23   34
                                         >> sigma = std( A )
                                         sigma =
                                         3.3040  5.1962  5.1962  3.5000

```

 WE MAKE LIVES BETTER  
**UT HEALTH SCIENCE CENTER**  
SAN ANTONIO



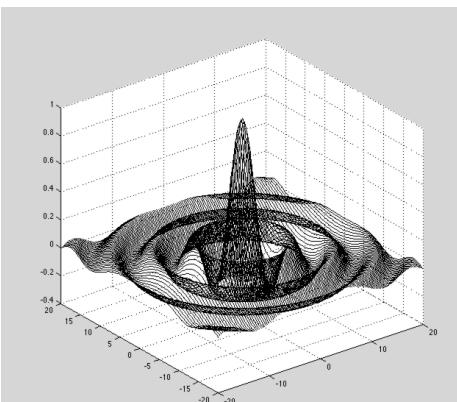
## Fancy plots

```

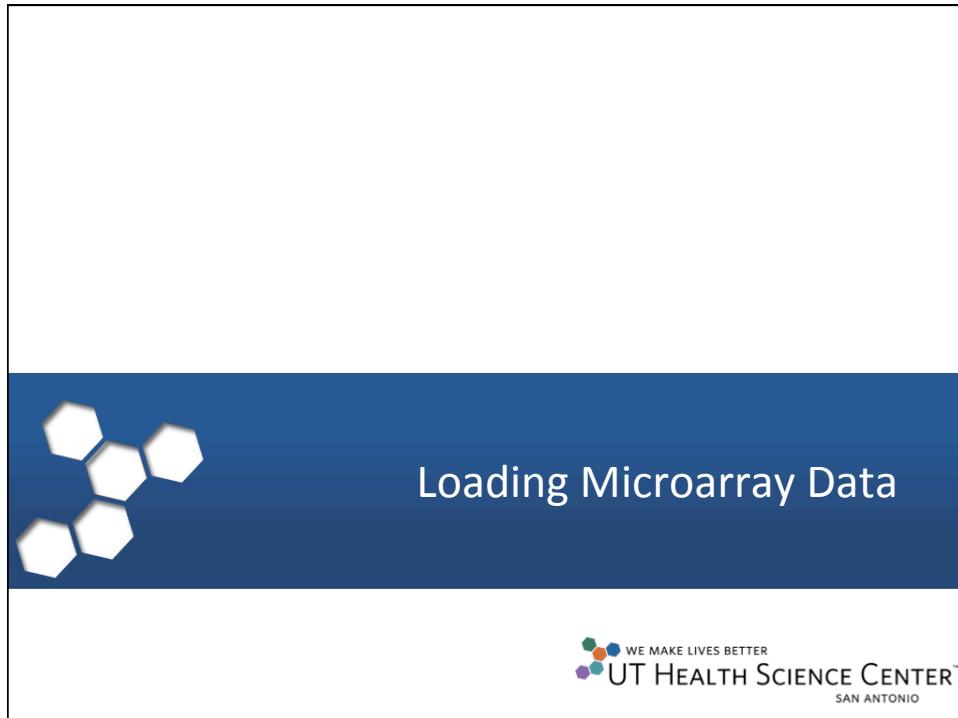
x1 = -20:0.3:20;
y1 = -20:0.3:20;
[x, y] = meshgrid(x1,y1);
r = sqrt( x.^2 + y.^2 );
s = sin(r)./r;
s(find(r==0)) = 1;
plot3( x, y, s );
grid on;

Try these:
plot3( x, y, s, 'k' );
mesh( s );
surf( s );
shading interp

```



 WE MAKE LIVES BETTER  
**UT HEALTH SCIENCE CENTER**  
SAN ANTONIO



### Reading generic text, tab-delimited file

- GSE21974\_AllGene.txt – Using excel.

	B	C	D	E	F
1	GenelD	GSM546381 GSM546382 GSM546383 GSM546384 GSM546385 G			
2	A_23_P100001 FAM174B	7.948285	10.721916	9.540191	10.191727
3	A_23_P100011 AP352	4.52518	3.7756774	4.857937	3.8732293
4	A_23_P100022 SV2B	4.905316	7.1086946	6.0053754	8.226337
5	A_23_P100056 RBPMS2	3.0645268	3.7501602	4.28707	5.0768857
6	A_23_P100074 AVEN	10.123294	8.855872	9.672125	9.154279
7	A_23_P100092 ZSCAN29	6.947331	6.3130016	7.057054	6.506809
8	A_23_P100103 VPS39	6.259719	6.1501107	6.065645	6.3619065
9	A_23_P100111 CHP	3.90169	3.8313572	3.7563396	3.9412627
10	A_23_P100127 CASC5	5.7190323	5.2640758	5.82081	3.444523
11	A_23_P100133 ATMIN	5.392776	5.567668	5.455074	4.897296
12	A_23_P100141 UNKL	7.5177383	8.090331	7.252972	8.595299
13	A_23_P100156 TMEM127	6.3529806	6.305977	6.958699	6.4668355
14	A_23_P100177 MMP15	4.5368176	2.1888537	3.286565	2.5291495
15	A_23_P100189 PRM1	3.2248104	1.7352371	3.4401698	2.3070922
16	A_23_P100196 USP10	9.556965	8.791511	9.282677	8.567565
17	A_23_P100203 HSBP1	9.32378	9.263293	9.220312	8.570824
18	A_23_P100220 ESRP2	10.165815	10.114584	10.044081	9.091385

Normal View Ready Sum=4.52518

IS BETTER  
UT HEALTH SCIENCE CENTER™  
SAN ANTONIO



## Importdata, load, textscan, fscanf, xlsread

```
importdata( 'GSE21974_AllGene.txt' );
```

Problem: because the first line is text, this command does not work

```
load 'GSE21974_AllGene.txt'
```

Use textscan:

```
fid = fopen('GSE21974_AllGene.txt' , 'r' );
A = textscan( fid, '%s', 58 );
fclose( fid );
```

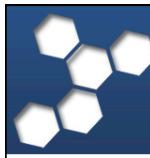
 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



## Read Microarray data

- Affyread: affymetrix CEL file
- Agferead: agilent feature extraction
- Ilmnbsread: illumina beadStudio output file (text file)

 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



## We use custom code

- For generic text
  - `ReadTextWithoutHeader()`
  - e.g.

```
data = ReadTextWithHeader( 'GSE21974_AllGene.txt', 3 );
```
- For agilent microarray
  - `ReadAgilentResult()`
  - e.g.

```
data = ReadAgilentResult(
'GSM546381_US91803681_251485049822_S01_GE1_105_Dec08_1_1.txt' );
```

**Request source code @ CBBI**





## Examination of Raw Agilent Microarray Data

- Download one of example Agilent gene expression data
  - <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM546381>

TYPE	FEPARAMS	DATA	TYPE	STATS	DATA	TYPE	FEATURES
text	Protocol_Name	GE1_105_Design(Read Only)	float	gDataOffsetAvg	6.16137	integer	FeatureNum
text	Protocol_date	10/17/09 12:35	float	gDataOffsetMedian	6.16137	integer	Row
text	Scan_Date	9/16/09 11:06	float	gDataOffsetStdDev	1.39238	integer	Col
text	Scan_ScannerName	Agilent Technologies Scanner G2505C	integer	gDataOffsetSumPnts		text	acessions
integer	Scan_NumChannels	1	integer	gSaturationValue	779631	text	chr_coord
float	Scan_MicronsPerPixelX	5	float	gAvgSig2Bkgd	589.932	integer	SubTypeMask
float	Scan_MicronsPerPixelY	5	float	gAvgSig2BkgdNegCtrl	1.47761	text	SubTypeName
text	Scan_FileGUID	d7d9523-fb3-43d3-bc16-4c98e5173	float	gAvgSig2BkgdQC_NegCtrl	399.247	integer	SubType
text	Grid_Name	014850_0_F_20090416	integer	gNumMisFeat		text	Sequence
text	Grid_Date	4/16/09 09.00	float	gLocalBGrillesNetAve	12.7466	integer	ProbeID
integer	Grid_NumSubGridRows	1	float	gLocalBGrillesAve	18.908	integer	ControlType
integer	Grid_NumSubGridCols	1	float	gLocalBGrillesDev	1.73806	text	ProbeName
integer	Grid_OffsetX	532	integer	gLocalBGrillesInfol	44645	text	Genotype
integer	Grid_OffsetY	85	float	gGlobalBGrillesAve	18.908	text	SystematicName
float	Grid_RowSpacing	73.3235	float	gGlobalBGrillesDev	1.73806	text	Description
float	Grid_ColSpacing	63.5	integer	gGlobalBGrillesNum	44645	float	PositionX
float	Grid_OffsetX	0	integer	gNumMisFeatInfol	1	float	PositionY
float	Grid_OffsetY	36.6617	integer	gNumPopnBGOL	138	float	gSampled
float	Grid_NonSpotWidth	65	integer	gNumPopnBGOL	0	boolean	gSampled
float	Grid_NonSpotHeight	65	integer	gNumPopnBGOL	370	float	gProcessedSignal
text	Grid_GenomicBuild	hg18-NCBI36-Mar2006	float	gOffsetUsed	6.16137	float	gProcessedSigError
text	FeatureExtractor_Barcodes	251485049822_1_1	float	gGlobalFeatInlierAve	2228.19	integer	gNumPixOLH
text	FeatureExtractor_Sample		float	gGlobalFeatInlierDev	11297.8	integer	gNumPixOLo
text	FeatureExtractor_ScanName	D:\Data\014850_0_F_20090416.xml	float	gGlobalFeatInlierNum	48876	integer	gNumPixPh
text	FeatureExtractor_ArrayName	US91803681_251485049822_S01	float	AnyColorPrcntFeatNonUnifOL	0.00222148	float	gMeanSignal
text	FeatureExtractor_DesignFileName	014850_0_F_20090416.xml	float	AnyColorPrcntBGNonUnifOL	0.00222140	float	gMedianSignal
text	FeatureExtractor_PrintingFileName		float	AnyColorPrcntFeatPopnOL	0.306564	float	gNxDev
text	FeatureExtractor_ExtractionTime	Agilent-014850	float	AnyColorPrcntFeatPopnOL	0.821868	float	gNxDevIDR
text	FeatureExtractor_ExtractionTime	9/16/09 12.07	float	TotalPrcntFeatOL	0.00239542	integer	gNcNmPx
text	FeatureExtractor_UserName	sa_agilent	integer	gNumNegFeatBG	4210	float	gNmSignal
text	FeatureExtractor_ComputerName	AGILENT02	integer	gNmCrNmNegFeatBGsubS	3344	float	gNmSignal
text	FeatureExtractor_ScanFileGUID	957a1497-1aa5-4f45-ba2e-2ab8537d	float	gSpatialDeterrMMSFit	1.17339	float	gNpNxDev

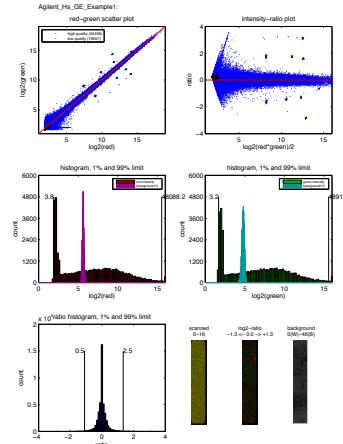


 FirstLook

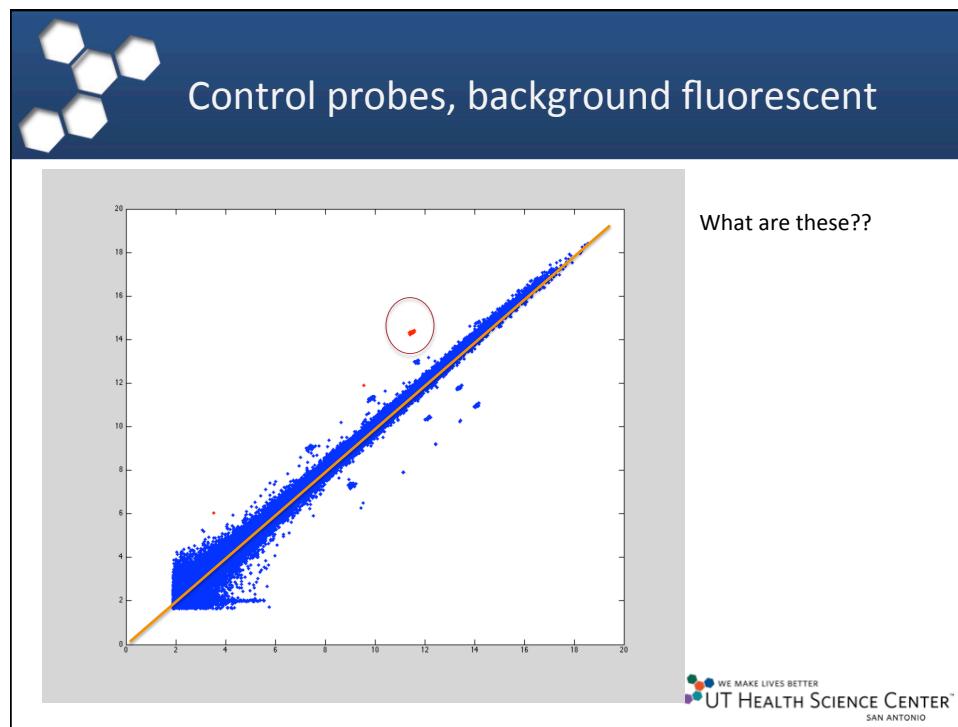
When you use  
`ReadAgilentResult('Agilent_Hs_GE_Example1.txt'),`

it generate a pdf file,  
`Agilent_Hs_GE_Example1.txt_qFig.pdf`

It is mainly for quality control purpose.



WE MAKE LIVES BETTER  
**UT HEALTH SCIENCE CENTER**  
 SAN ANTONIO





## Normalization

- Boxplot, Lowess, quantile

```
>> boxplot( data.value, 'plotstyle', 'compact' );
>> data.normValue = quantilenorm( data.value );
```

(for agilent microarrays, lowess correction is not needed)

```
hold on;
plot( data.value(8698,:), 'r' );
text( 1, data.value(8698,1), data.GeneID(8698), 'interpret',
'none', 'horizontalalignment', 'right');
```

 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



## Intensity, median/mean, ratio, etc

- This is single channel data file with data in `gProcessedSignal`
- For 2 channels, data are in `rProcessedSignal`, and `gProcessedSignal`. With ratio in log10 transformation.
  - <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM785046>

 WE MAKE LIVES BETTER  
UT HEALTH SCIENCE CENTER  
SAN ANTONIO



## PCA, MDS, etc

- Principal component analysis  
`[coefs, scores, sigma2] = princomp( data.value' );`
- Multidimensional scaling



## Differential Gene Expression

```
p = mattest( data.value( 1:2:57 ), data.value(2:2:57) );
adjp = mafdr( p, 'BHFDR', true );

Or

[fdr, q] = mafdr( p );
```



