

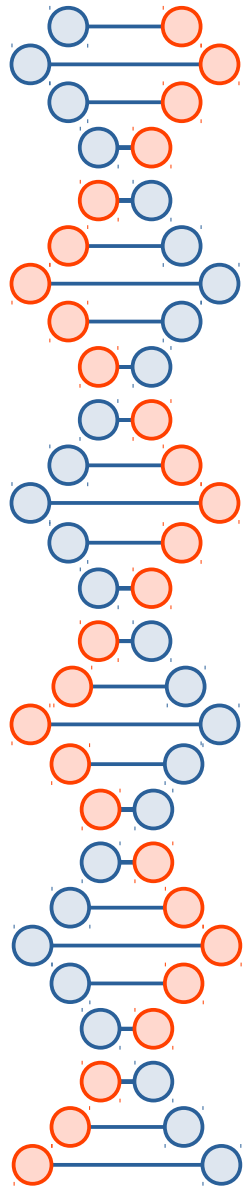


# Omics Data Analysis

(Lunch & Learn)

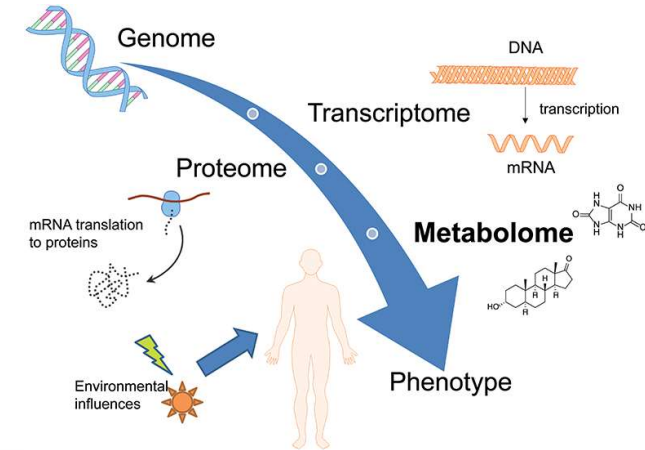
Kith Pradhan

Nov 17, 2025

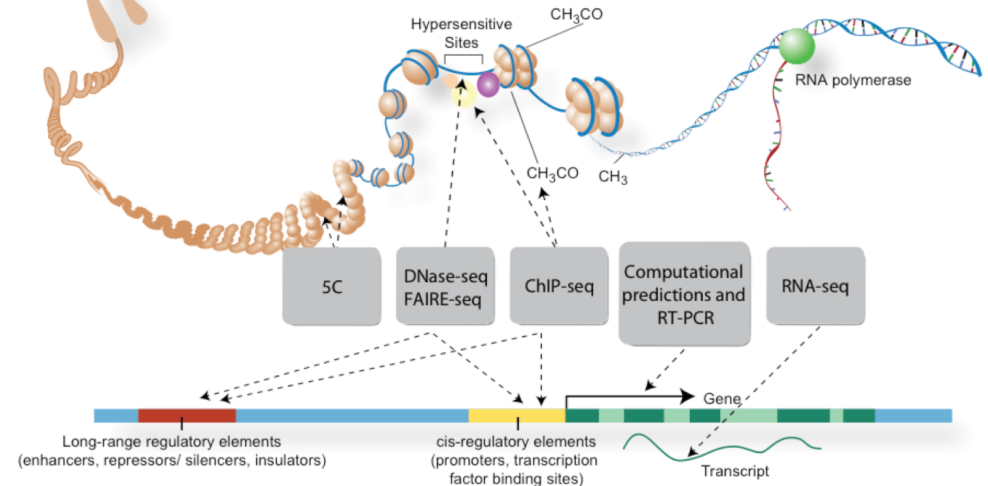


# What is Omics

- The study of biological molecules at large scale
  - Genomics
  - Transcriptomics
  - Proteomics
  - Metabolomics
  - Epigenomics
  - Lipidomics
  - Microbiomics



<https://www.creative-proteomics.com/resource/what-is-metabolomics.html>

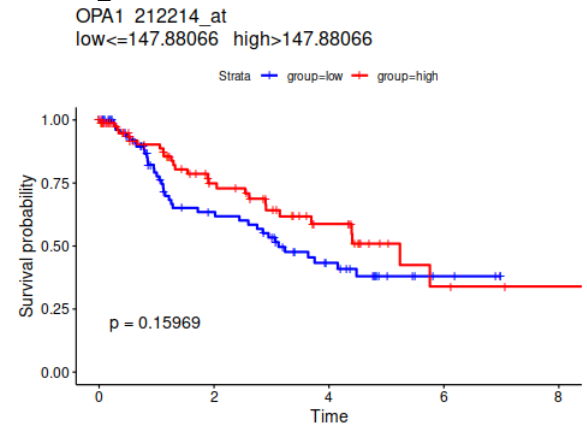
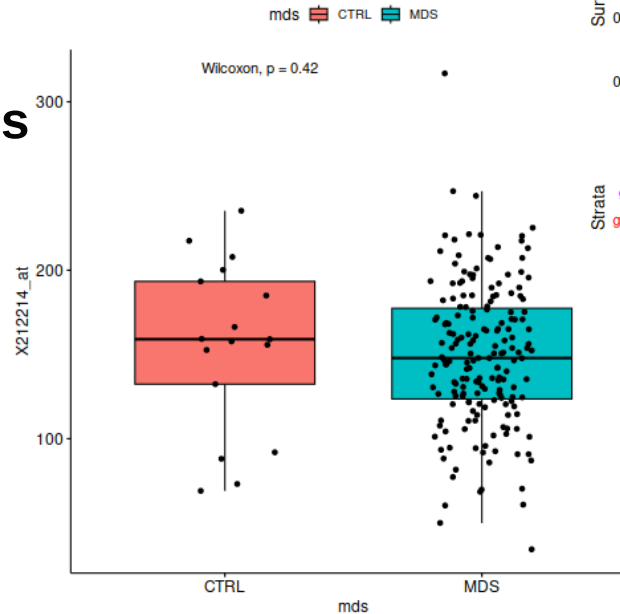


**ENCODE Project**

# Example: Microarray

- The study of biological molecules at large scale

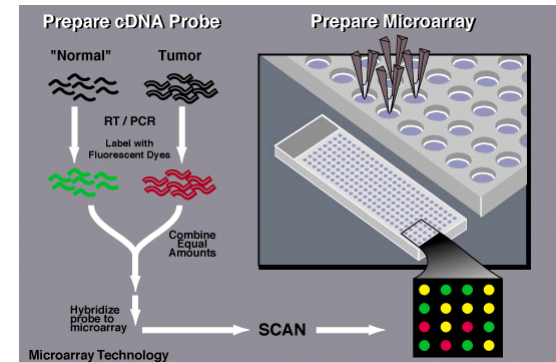
- Genomics
- **Transcriptomics**
- Proteomics
- Metabolomics
- Epigenomics
- Lipidomics
- Microbiomics



Number at risk

Strata	group=low	88	38	18	4	0
	group=high	88	38	18	4	1

Time

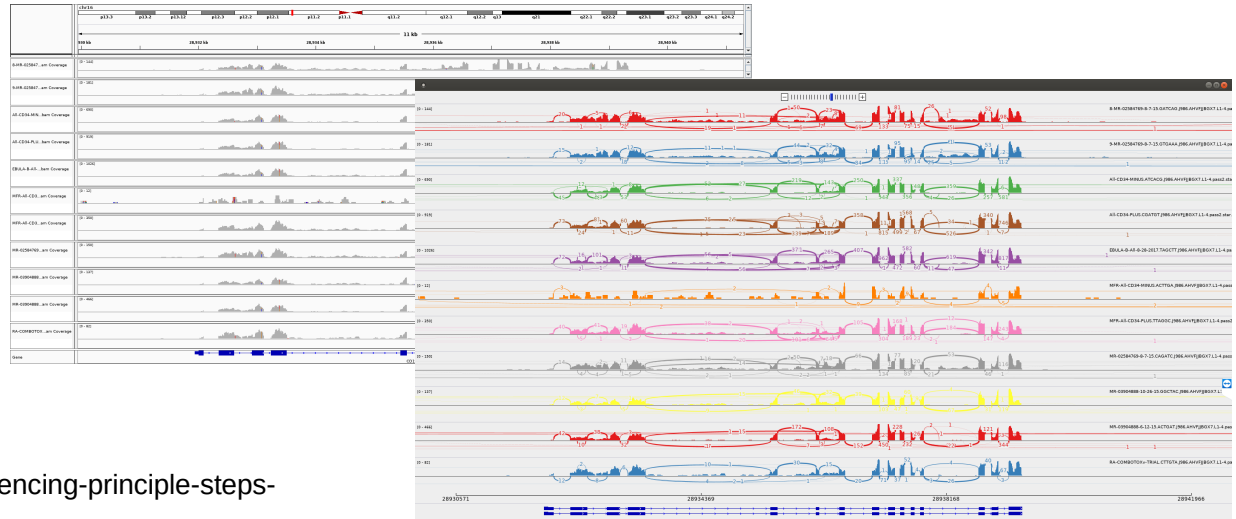
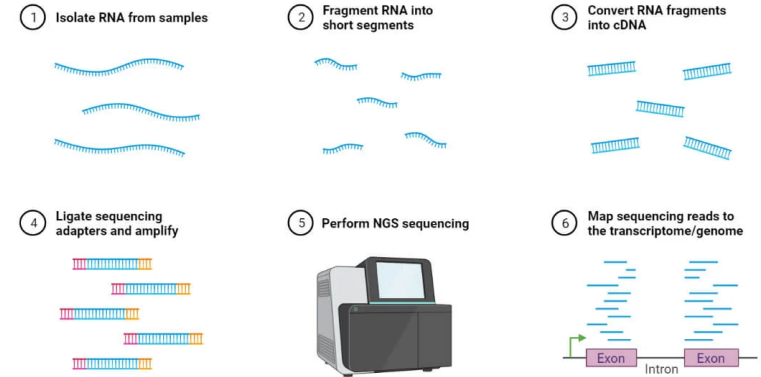


# Example: NGS RNA-seq

- The study of biological molecules at large scale

- Genomics
- **Transcriptomics**
- Proteomics
- Metabolomics
- Epigenomics
- Lipidomics
- Microbiomics

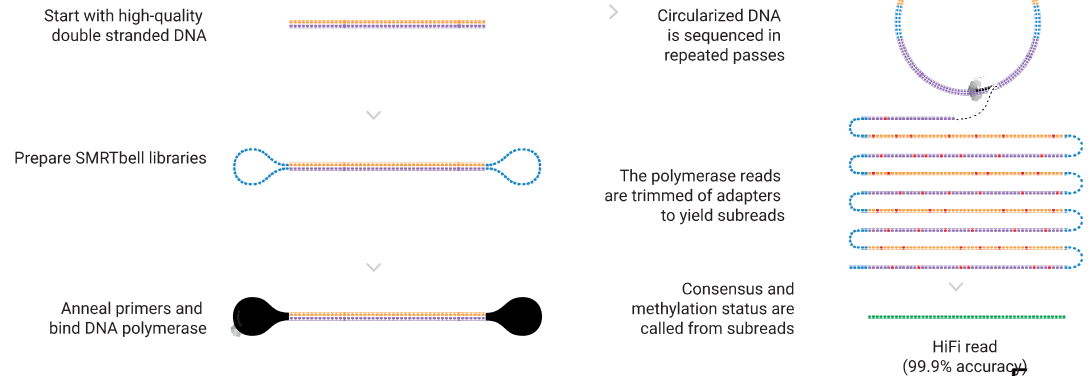
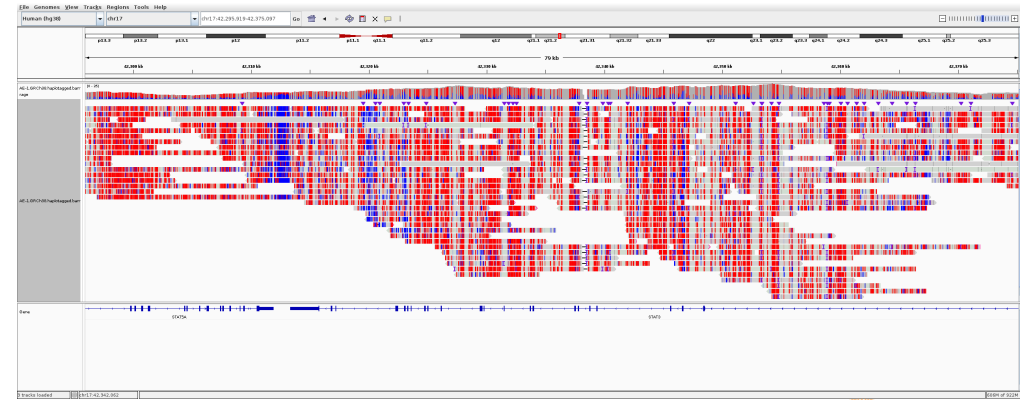
## RNA Sequencing



<https://microbenotes.com/rna-sequencing-principle-steps-types-uses/>

# Example: PacBio HiFi-seq

- The study of biological molecules at large scale
  - Genomics
  - **Transcriptomics**
  - Proteomics
  - Metabolomics
  - **Epigenomics**
  - Lipidomics
  - Microbiomics

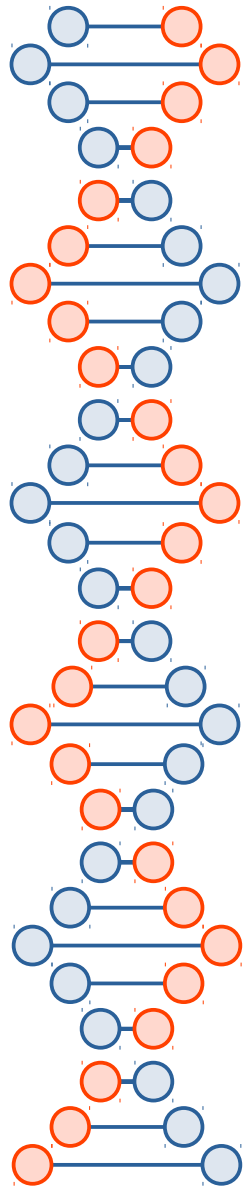


<https://www.pacb.com/technology/hifi-sequencing/>

# The Biggest Issue in Omics Analysis?

- The differences between Omics technologies are vast and each requires specialized tools to process and analyze.
- But there is a common issue shared across all.
  - # of variables: N
  - # of observations: M
  - $N \gg M$

ENSEMBL	ENTREZGENE	SYMBOL	Mock1.1	Mock2.1	ZIKV1.1	ZIKV2.1	Mock1.2	Mock2.2	ZIKV1.2
ENSG00000010030	51513	ETV7	0.0708594	0	0.3114447	0.2252182	0.0493682	0.0367962	0.2537942
ENSG00000010072	83932	SPRTN	3.1900906	2.6378208	3.0984427	3.2892438	3.6348061	3.2145277	4.6315328
ENSG00000010165	51603	METTL13	4.9670487	4.2128095	3.6224677	3.5393062	4.4813932	4.3178785	3.6262744
ENSG00000010219	8798	DYRK4	2.0031906	2.0030852	3.4851214	2.9813266	2.0685315	2.0866457	3.0435702
ENSG00000010244	7756	ZNF207	12.419449	12.024094	12.350847	13.169165	15.314333	14.552613	16.34819
ENSG00000010256	7384	UQCRC1	40.287622	40.417455	40.098703	38.056824	38.66833	40.203852	43.152384
ENSG00000010270	83930	STARD3	19.151117	17.581794	20.533327	19.142134	22.498407	21.205717	25.820423
ENSG00000010278	928	CD9	1.3767652	1.0982165	0.3922089	0.7563255	0.9651221	1.1430112	0.4565833
ENSG00000010282	57467	HHATL	0	0.0626603	0	0	0.0202698	0.020144	0.0208408
ENSG00000010292	9918	NCAPD2	32.54455	29.572793	14.275639	14.808391	30.03627	28.103656	14.74612
ENSG00000010295	25900	IFFO1	4.9790368	6.1973561	6.5368079	6.3027053	5.0757457	5.4397063	6.9248392
ENSG00000010310	2696	GIPR	0.4262339	0.3824975	0.3122339	0.3386833	0.4536887	0.3842655	0.368934
ENSG00000010318	51533	PHF7	2.5745529	2.5632081	1.3344099	1.6468745	2.3408447	2.3207032	1.1192984
ENSG00000010319	56920	SEMA3G	0.0671404	0.1446025	0	0	0.0779619	0.0958788	0.0721422
ENSG00000010322	11188	NISCH	29.013	32.390026	41.020053	39.998457	25.416185	27.365164	38.895627
ENSG00000010327	23166	STAB1	0	0	0.0174577	0.0168325	0.0147589	0.0068753	0.032246
ENSG00000010361	80199	FUZ	5.6453597	6.5091477	3.4462152	3.7456973	4.8336793	5.6212736	3.635468
ENSG00000010379	6540	SLC6A13	0.0792931	0.042694	0.2613846	0.4200393	0.082866	0.0892141	0.3408008
ENSG00000010404	3423	IDS	16.637373	17.058433	19.773289	21.232514	18.890098	17.876043	25.156345
ENSG00000010438	5646	PRSS3	0.2021879	0.1088647	0	0.321315	0.0939104	0.0874943	0.1689724
ENSG00000010539	7752	ZNF200	1.8396994	1.8361512	1.9475265	2.020396	2.1310069	2.0951206	2.5794192
ENSG00000010610	920	CD4	0.2936423	0.1976334	0.201861	0.1166633	0.2727766	0.3017913	0.184052
ENSG00000010626	10233	LRRRC23	3.8557374	3.8555346	3.08676	2.2759283	3.1532244	3.4991297	2.8737798
ENSG00000010671	695	BTK	0	0	0.0404027	0	0.0256175	0.0318231	0.0219493
ENSG00000010704	3077	HEF	0.0502817	0	0.0651427	0.0628088	0.013768	0.0128274	0.0380285



$$N \gg M$$

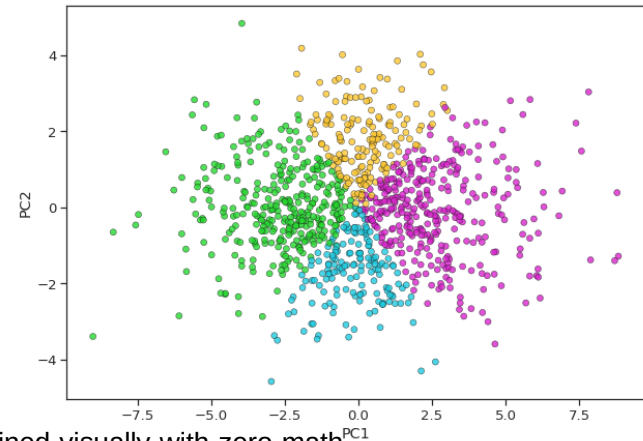
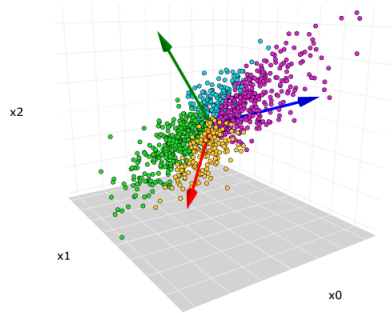
- How do we address this?
1. Dimensionality reduction
  2. Choosing the right test
  3. Dealing with the multiple comparisons in N variables

ENSEMBL	ENTREZGENE	SYMBOL	Mock1.1	Mock2.1	ZIKV1.1	ZIKV2.1	Mock1.2	Mock2.2	ZIKV1.2
ENSG00000010030	51513	ETV7	0.0708594	0	0.3114447	0.2252182	0.0493682	0.0367962	0.2537942
ENSG00000010072	83932	SPRTN	3.1900906	2.6378208	3.0984427	3.2892438	3.6348061	3.2145277	4.6315328
ENSG00000010165	51603	METTL13	4.9670487	4.2128095	3.6224677	3.5393062	4.4813932	4.3178785	3.6262744
ENSG00000010219	8798	DYRK4	2.0031906	2.0030852	3.4851214	2.9813266	2.0685315	2.0866457	3.0435702
ENSG00000010244	7756	ZNF207	12.419449	12.024094	12.350847	13.169165	15.314333	14.552613	16.34819
ENSG00000010256	7384	UQCRC1	40.287622	40.417455	40.098703	38.056824	38.66833	40.203852	43.152384
ENSG00000010270	83930	STARD3	19.151117	17.581794	20.533327	19.142134	22.498407	21.205717	25.820423
ENSG00000010278	928	CD9	1.3767652	1.0982165	0.3922089	0.7563255	0.9651221	1.1430112	0.4565833
ENSG00000010282	57467	HHATL	0	0.0626603	0	0	0.0202698	0.020144	0.0208408
ENSG00000010292	9918	NCAPD2	32.54455	29.572793	14.275639	14.808391	30.03627	28.103656	14.74612
ENSG00000010295	25900	IFFO1	4.9790368	6.1973561	6.5368079	6.3027053	5.0757457	5.4397063	6.9248392
ENSG00000010310	2696	GIPR	0.4262339	0.3824975	0.3122339	0.3386833	0.4536887	0.3842655	0.368934
ENSG00000010318	51533	PHF7	2.5745529	2.5632081	1.3344099	1.6468745	2.3408447	2.3207032	1.1192984
ENSG00000010319	56920	SEMA3G	0.0671404	0.1446025	0	0	0.0779619	0.0958788	0.0721422
ENSG00000010322	11188	NISCH	29.013	32.390026	41.020053	39.998457	25.416185	27.365164	38.895627
ENSG00000010327	23166	STAB1	0	0.0174577	0.0168325	0.0147589	0.0068753	0.032246	0
ENSG00000010361	80199	FUZ	5.6453597	6.5091477	3.4462152	3.7456973	4.8336793	5.6212736	3.635468
ENSG00000010379	6540	SLC6A13	0.0792931	0.042694	0.2613846	0.4200393	0.082866	0.0892141	0.3408008
ENSG00000010404	3423	IDS	16.637373	17.058433	19.773289	21.232514	18.890098	17.876043	25.156345
ENSG00000010438	5646	PRSS3	0.2021879	0.1088647	0	0.321315	0.0939104	0.0874943	0.1689724
ENSG00000010539	7752	ZNF200	1.8396994	1.8361512	1.9475265	2.020396	2.1310069	2.0951206	2.5794192
ENSG00000010610	920	CD4	0.2936423	0.1976334	0.201861	0.1166633	0.2727766	0.3017913	0.184052
ENSG00000010626	10233	LRRRC23	3.8557374	3.8555346	3.08676	2.2759283	3.1532244	3.4991297	2.8737798
ENSG00000010671	695	BTX	0	0	0.0404027	0	0.0256175	0.0318231	0.0219493
ENSG00000010704	3077	HEF	0.0502847	0	0.0651427	0.0628088	0.013768	0.0128274	0.0380285



# 1. Dimensionality Reduction

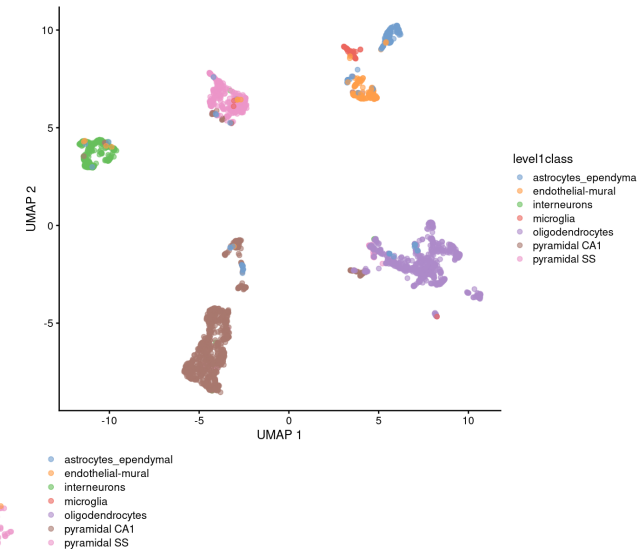
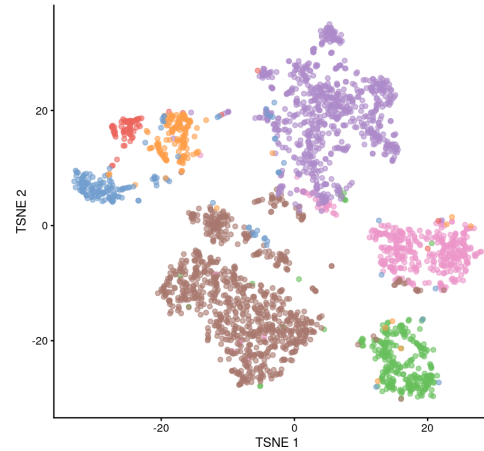
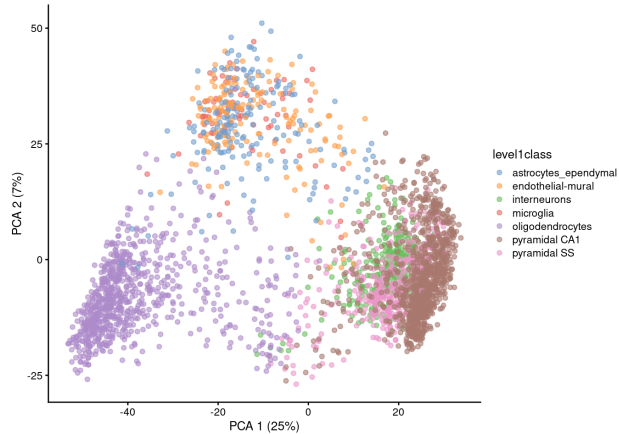
- Huge # of variables to look over
- Hard to show more than 3 vars with traditional plots
- Not all the variables are equally important
- Can we get away with a smaller representation that only contains relevant informatio





# Dimensionality Reduction

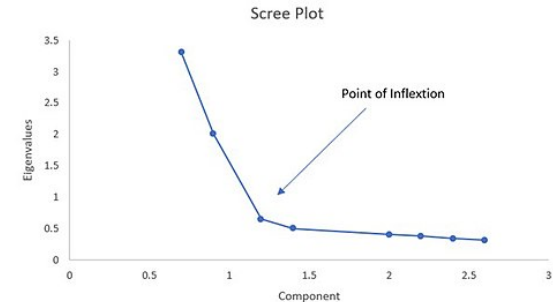
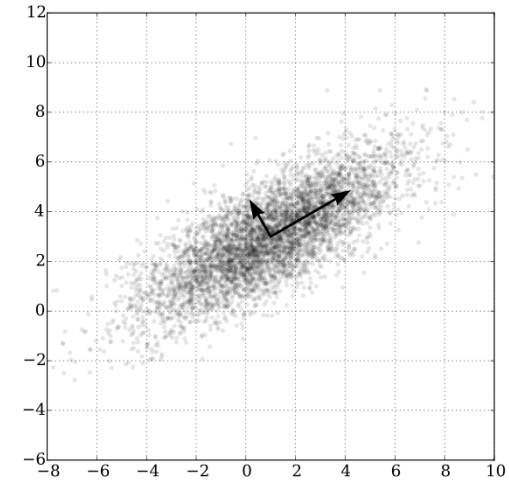
- Common dimensionality reduction techniques:
  - PCA
  - t-SNE
  - UMAP

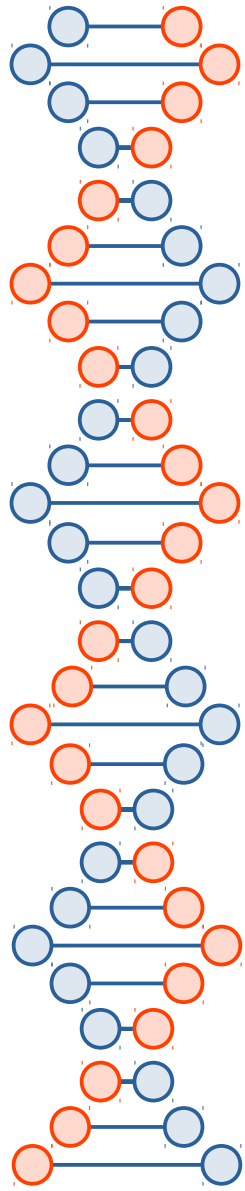


# Principle Component Analysis(PCA)

(Pearson 1901)

- Linear projections along lines of greatest variance
- Principal components(the new orthogonal axes) are ordered by the amount of variance explained.
- The idea is that we can use a smaller set of these PCs to capture most of the information of the original data.

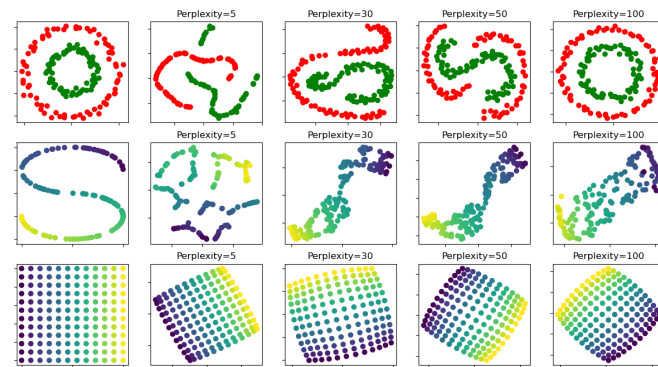


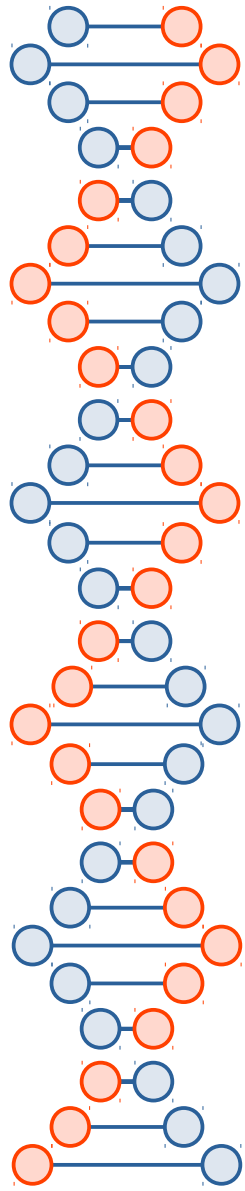


# t-distributed Stochastic Neighbor Embedding (t-SNE)

(Maaten & Hinton 2008)

- A Non-linear approach to dimension reduction
- Idea: Distances between the points in the high dim space should be similar to distance between points in the low dim space
  - High dim: Distances represented by Gaussian
  - Low dim: Distance represented by t-distribution
  - Loss function: Measures differences between the distance representations in high and low dimensions
    - Kullback-Leibler divergence
- Optimize: find a low dim representation that minimizes the loss function.
  - Start with a random low dim representation
  - Iteratively shift points in the low dim space so that loss function goes down
    - move along the gradient of the Kullback-Leibler divergence
- Key hyperparameter
  - Perplexity: The width of the gaussian or t-distribution.
    - A wider bell shaped curve means we're looking at more points at each step.

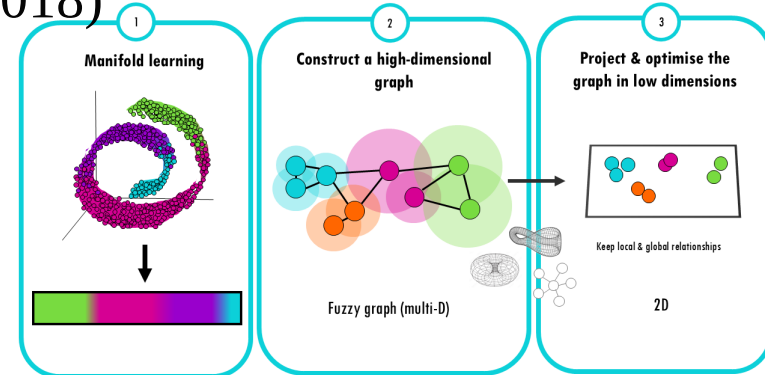




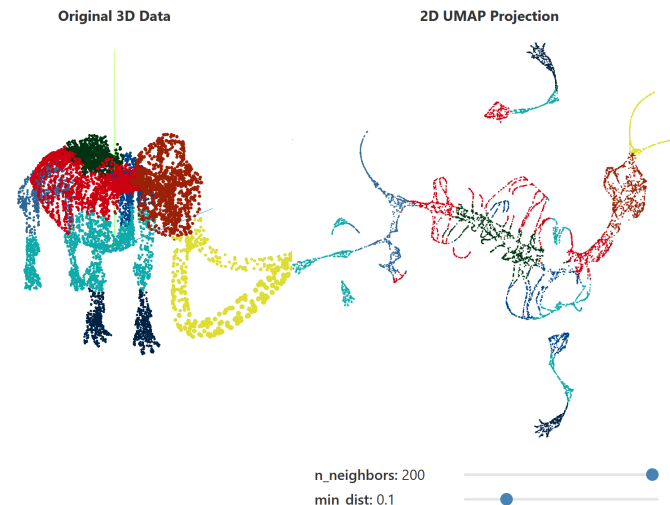
# Uniform Manifold Approximation and Projection (UMAP)

(Innes & Healy 2018) How does UMAP summarise many dimensions into 2?

- A Non-linear approach to dimension reduction
- Idea: Distances between the points in the high dim space should be similar to distance between points in the low dim space
  - High dim: weighted sum of k nearest neighbor graphs
  - Low dim: weighted sum of k nearest neighbor graphs
  - Loss function: Measures differences between the distance representations in high and low dimensions
    - Cross entropy of adjacency matrices
- Optimize: find a low dim representation that minimizes the loss function.
  - Start with a random low dim representation
  - Iteratively shift points in the low dim space so that loss function goes down
    - Minimize cross entropy by stochastic gradient descent
- Key hyperparameter
  - n\_neighbors: how many neighbors to look at



<https://biostatsquid.com/umap-simply-explained/>



<https://pair-code.github.io/understanding-umap/>

# Dimensionality Reduction

- Pros

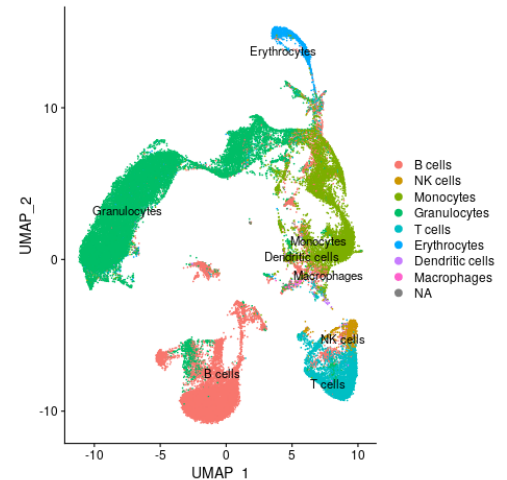
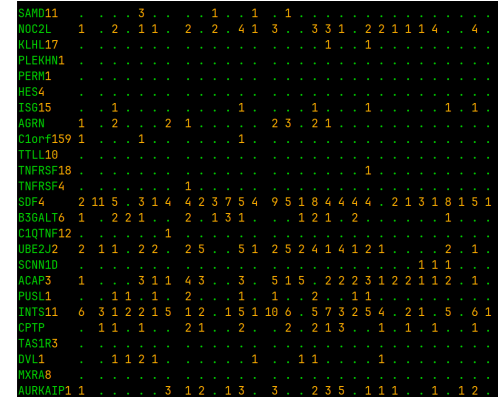
- PCA
  - Easy to interpret
  - Fast
- T-SNE
  - Captures non-linear relationships between variables
  - Preserves local structure
- UMAP
  - Captures non-linear relationships between variables
  - Preserves local and global structure
  - Scalable
  - Stable

- Cons

- PCA
  - Only linear relationships
- t-SNE
  - Very slow
  - Changing the perplexity hyperparameter gives very different results across runs
- UMAP
  - Not as good as t-SNE if you're ONLY interested in local structure.

# Dimensionality Reduction: example1

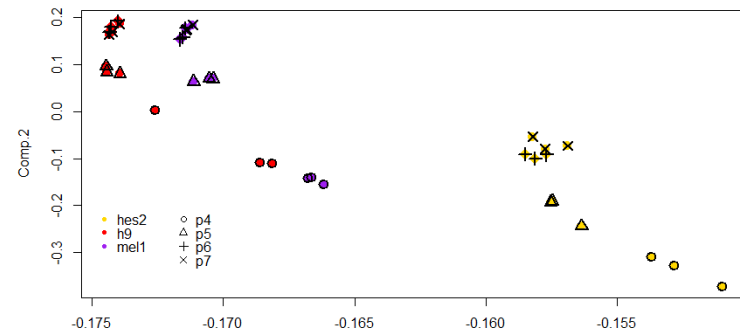
- Single-cell sequencing
  - We are able to capture gene expression for each individual cell
  - Too many genes to be able to map out everything
  - A 2D UMAP (or t-SNE) of all cells allows us to:
    - Visually examine the individual cells
    - Pick out clusters of similar cells
    - Identify celltypes of the clusters and their relationships to one another
      - Manual celltype type ID
      - Automatic ID by Machine learning algorithms



# Dimensionality Reduction: example 2

- RNA-seq experiment with a batch effect
- Examining pluripotency of 3 different cell lines
- A PCA plot of the samples gives us a global picture of gene expression
- This allows us to identify batch effects visually
- Once you know they are there, they can be corrected.

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Annotation	hes2.p4.r1	hes2.p4.r2	hes2.p4.r3	hes2	hes2.p4.r2	hes2.p4.r3	hes2	hes2.p4.r2	hes2.p4.r3	hes2.p7.r1	hes2.p7.r2	hes2.p7.r3
2	CellLine	hes2	hes2	hes2	hes2	hes2	hes2	hes2	hes2	hes2	hes2	hes2	hes2
3	Fraction	p4	p4	p4	p5	p5	p5	p5	p5	p5	p7	p7	p7
4	RAM_1004603	11.0095455876	10.4599950613474	10.8365305689498	10.032987008184	9.52434373400795	9.813024771933504	8.80054962195996	9.04217235317493	8.97661185447159	8.73165944875181	8.72752343193379	8.90548729331
5	RAM_1001799	13.743094660574	13.599187463306	13.920307223828	13.046466957404	13.024347117111	13.795899048464	13.905017573738	13.448070307523	13.70535793179	13.748263680808	13.699434399751	13.676610352
6	RAM_1827636	8.33732456171098	8.23112115678901	8.50260608181731	8.303203906837	8.296239917535	8.1746106574821	8.19600779535305	8.6293876705453	8.53405953976998	8.73369032506631	8.22601071928101	8.49034950275
7	RAM_1812262	12.5211014823021	12.2093126866689	11.8294320282968	11.629625011298	11.8987349381152	11.9471122295206	11.873030710344	11.8405725798907	11.8266213087159	11.792747119844	11.693727233214	11.41053089512
8	RAM_1712803	12.9884667298018	12.9510148878381	12.7151605285208	12.8741460778703	12.8779799378764	13.004744230155	12.8424953840332	12.8720101790974	12.829678310747	12.796572003172	12.747269507787	12.8684850913
9	RAM_1815941	10.4461788746742	10.2528329487879	10.074452348191	10.771953622004	10.7362791972383	10.3412582446824	9.970990930665	10.470024954589	10.3289638246542	11.0891054650974	10.6736448276596	10.425338308
10	RAM_1786374	8.693531847296	8.2893367072029	8.6104325642491	8.602767159495	8.110047744026	8.154545499149	7.87829761067878	8.234347314994	8.0321462189708	7.7533802569201	8.0186225536526	7.692141705
11	RAM_1674380	10.6164349749513	10.4893203248763	10.721347988159	10.6327456572347	10.430205970987	10.541338679708	10.2548266999451	10.361624699451	10.500742311246	10.186838432399	10.41396236627	10.303524558
12	RAM_1669562	11.047138639553	11.19777467918	11.298735630111	10.4548057197683	10.905703083117	10.955995481516	10.675752939001	10.75248470484	10.8989890738	10.63726876461	11.026959878	
13	RAM_1849494	10.65862923458	10.588213437249	10.765397725984	10.34831184854	10.56522020642	10.464873716189	10.254269696921	10.442697091512	10.582959489707	10.143137434797	10.44320277287	10.562497205
14	RAM_1785324	10.2152397892424	10.67734321003	10.373649117515	10.5964772002497	10.7018623484346	10.5691435765451	11.388313334432	11.064714700483	11.038541096338	11.1338282002178	11.29950854805	11.105901311
15	RAM_1786429	9.3527299595704	8.8648425425671	9.5801357765987	8.9405189950076	8.891978757828	8.9702277800787	8.9234959721695	9.2056448299312	9.3143785263543	9.348455209644	9.47030451668	9.366449567
16	RAM_1808821	10.0915981374973	10.7486422973431	10.622661360005	10.731813988267	10.6197023531	10.480324399444	10.766411577892	10.679302824571	10.7118071159893	10.867098408749	10.842140784531	10.927091595
17	RAM_1815973	7.731335254624	7.6784877048224	7.6737349619672	7.860771653135	8.0074708235297	7.799395624264	8.2211287009706	8.16826149289105	8.081763399962	8.05612319705188	8.2509448726651	8.23896241217
18	RAM_1798579	10.3832738761454	10.378693066273	10.278677352056	10.180412913771	10.128984447781	9.874789387414	10.22321787667	10.444011118989	10.466289677376	10.45950827576	10.166066175	
19	RAM_1718769	12.2371760819318	12.159537632628	12.87849184763	12.031968216109	12.905036188672	12.4814441789239	12.195615712134	12.2084248961375	12.409437456177	11.9827321484252	12.1724405311806	12.250764731
20	RAM_1827741	7.994034884945	8.0011725957953	7.978957272828	8.289034315708	8.498771440382	8.51032777639569	8.809393006262	8.455042698916	8.702529184663	8.7489621197144	8.686331847296	8.538769816
21	RAM_1876824	12.08246232357	12.3896195279791	12.127911053109	12.235798452468	12.336248541292	12.382729838342	12.190021980674	12.130353959503	11.6658839528328	11.827384395955	11.9296252011	
22	RAM_1739602	10.0647057784	10.22755779737	10.261195781666	10.224078334702	10.164841384281	10.506129676332	10.372837897234	10.3827650749815	10.28151729287	10.220237357465	10.302016938774	9.9552664652
23	RAM_1858677	11.4536244867693	11.476777770273	11.77139818636	11.025334751895	11.068140733959	11.259613861543	10.962032927778	10.820678633239	10.9962315434802	11.118303611194	10.9319938996164	10.8543364806
24	RAM_1801429	9.8163713085466	9.8395004288171	9.703817987927	9.5839394484327	9.800793327826	9.7421129648819	9.69102373661763	9.84071734814623	9.198760146395	9.810035048162	9.9156428928762	10.179601741
25	RAM_179609	12.261720181113	12.230484953496	12.028498171347	12.469884991784	12.295729416506	12.396878900876	12.50821089389	12.381824436306	12.327765187896	12.416786756562	12.29121428406	12.508901463
26	RAM_1801735	10.3886362154628	10.614393749513	10.310541781149	9.875105950104	9.8306111972953	10.472059595976	9.893036614623	10.507744313708	10.18797042108	10.135773404495	10.239616417	
27	RAM_1743635	8.9736977324788	8.694967234776	8.59618975614441	8.7807910924932	8.507584254204	8.5015386639475	8.6862324621637	8.484302830284	8.7689738212959	8.66018831718	8.8425637029471	8.5166212131
28	RAM_1704146	10.535787722573	8.9320589662271	8.334602499959	8.3802445947818	8.102248711871	8.3874511849161	8.4862111569318	8.9265012120678	8.3715218138538	8.8047791805167	8.478044884675	8.4331158412
29	RAM_1802843	10.8474347408484	11.085506547276	10.755987402081	10.75102216978	10.8774037957	10.693083975056	10.71193620188	10.793108778752	10.7811001145576	10.944339829628	10.9033296603	
30	RAM_1726059	10.659596754853	10.79652112078	10.786250125793	10.704773952013	10.74778931013	10.79767956469	11.19544259818	10.95049137848	10.97689295474	10.95104517601	10.7474646664	
31	RAM_1746532	8.3407464919886	9.05599405133893	9.9853038004585	8.262824324304	9.1309168183578	9.1189575864144	9.956622889393	9.1158895896516	8.973297312558	8.87325514768	9.0479077271596	9.56987627234
32	RAM_1750150	1.6820069589392	8.92531894992039	8.808240232594	9.3549584999993	9.106664747938	9.1350212415033	9.7619706069642	8.9418927403987	8.866107307743	8.8642910089437	9.00571027964	
33	RAM_1860732	11.840025350104	11.841394711196	11.89718962434	11.8995254096541	11.6040653786224	11.873303571777	11.846684164784	11.963235961294	11.8724680157533	11.8544923085437	11.800000000	
34	RAM_1728512	10.3737877589092	10.0842316705202	10.4634477408336	10.2459685798059	10.03311896483	8.8116289229184	10.2166246182237	10.155953513779	10.0378848669992	10.664185770596	10.688180173468	10.665339917
35	RAM_1796300	11.8228405343	11.8248327198532	11.8248792071874	11.913995751289	11.8854377434395	12.0255317964783	11.9071425972983	11.8449499287431	12.0157257443109	11.8233207949219	12.008349954893	12.0483790197



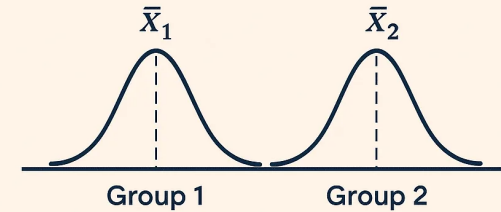
Comp.1



## 2. Choosing the right test

- A t-test is good for large sample sizes or if data is normally distributed
  - This is usually not the case in Omics datasets
- Each Omics technology needs its own particular methodology
  - 1. A stat that fits the data
  - 2. An "information borrowing" technique to deal with the  $N \gg M$  problem

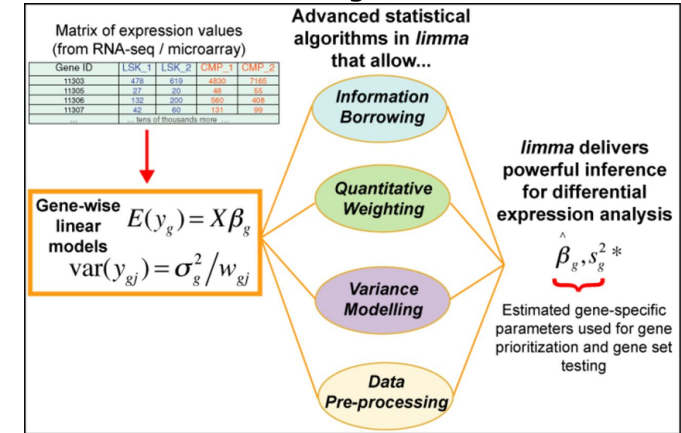
### Independent Sample t-Test



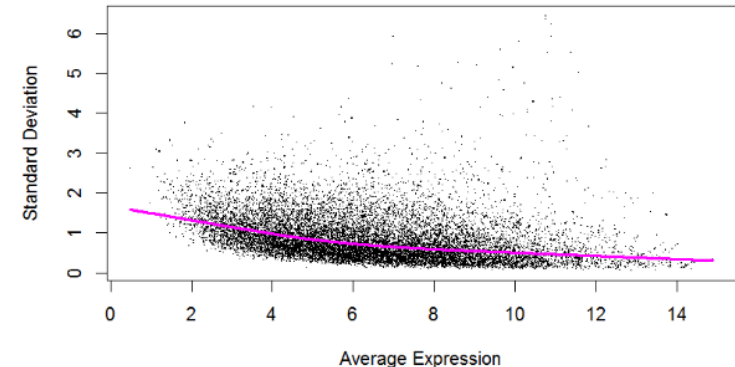
Hypotheses	Formula
$H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$	$t = \frac{\bar{X}_1 - \bar{X}_2}{SE}$
Standard Error	Pooled Variance
$SE = \sqrt{S_p \left( \frac{1}{n_1} \right)}$	$S_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$

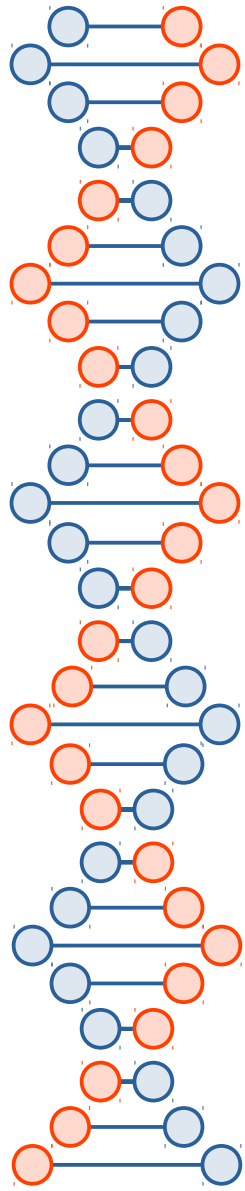
# Example 1: Linear Models for MicroArray data (LIMMA)

- The underlying statistical framework uses empirical Bayes linear modeling.
- Conceptually, for simple two group comparisons, it can be thought of as a moderated t-test
  - Looks for differences in gene expression across experimental groups.
  - But uses an empirical Bayesian approach to "shrink" variance



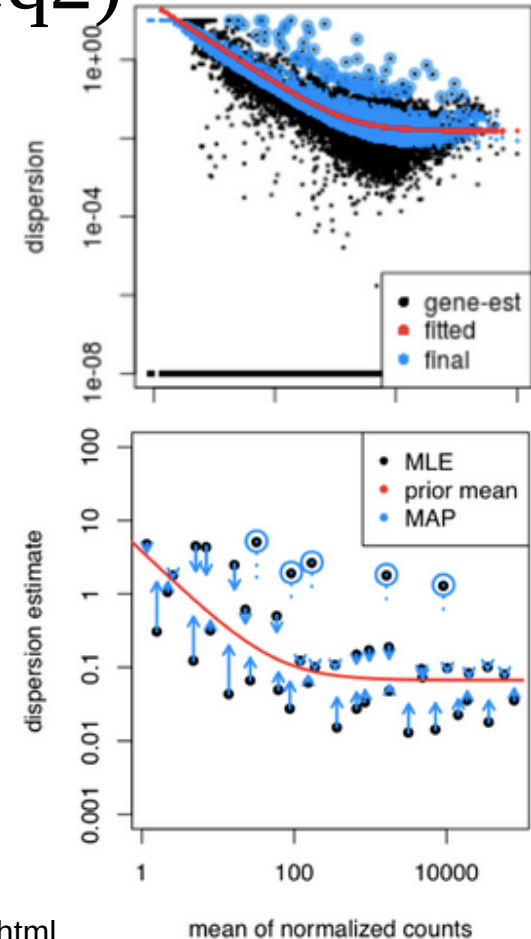
$$T_{\text{gene}} = \frac{\hat{D} - \hat{C}}{f(\text{Var}(\hat{D}, \hat{C}) + \alpha)}$$





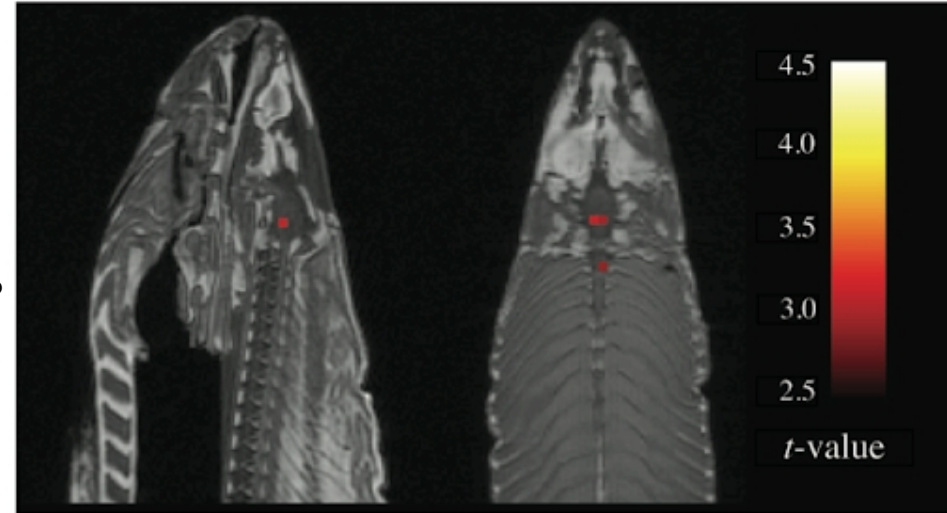
## Example 2: Differential Expression analysis for Sequencing data (DESeq2)

- RNA-seq gives us integer counts of the number of read fragments that map within a gene's boundary
- Negative Binomial Regression
  - Like Poisson, but allows for over/under dispersion through parameter ( $\theta$ )
  - Uses similar genes to estimate the dispersion parameter ( $\theta$ )

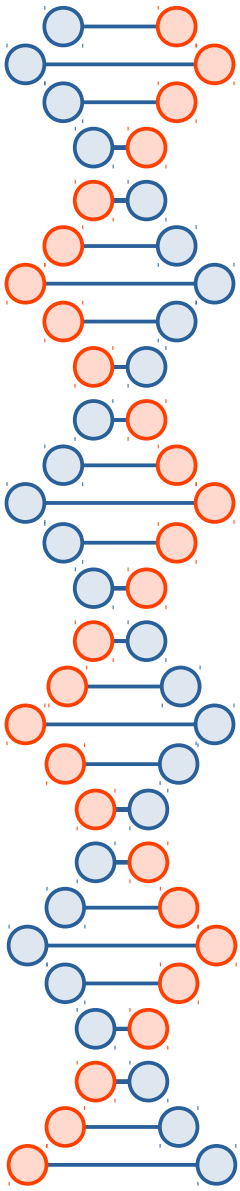


### 3. Dealing with multiple comparisons

- Type 1 error: Incorrectly rejecting the null hypothesis when the null is true.
  - Finding signal when there really is none there
- If you set significance level ( $\alpha$ ) at 0.05, and run 50,000 tests across all genes
  - How many false positives do you expect?
    - $50000 \times 0.05 = 2500$
- How do you address it?
  - Do nothing
  - Control for FWER
  - Control for FDR



<http://prefrontal.org/files/posters/Bennett-Salmon-2009.pdf>



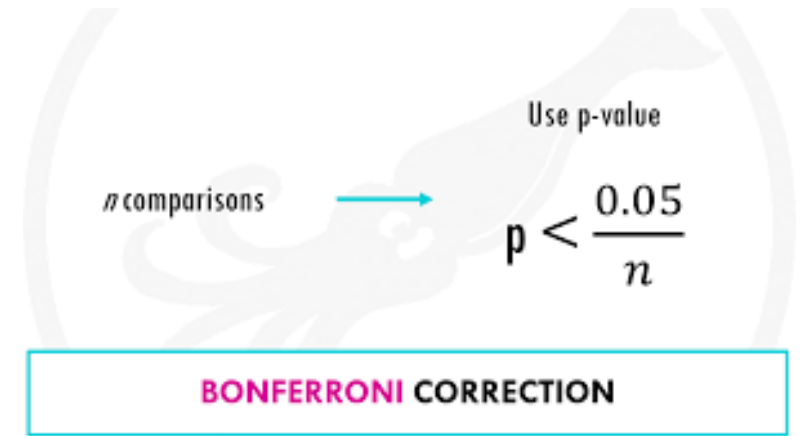
# Dealing with multiple comparisons

- Do nothing:
  - Ok if you have a specific location of interest known beforehand
    - Do NOT look at all the results then select the location of interest
  - Ok if there are a small set of interesting locations known beforehand
    - Show them all
  - Do NOT only show the significant results

# Dealing with multiple comparisons

## Control for Family-Wise Error Rate(FWER):

- Make at most X false positives across all tests
- How many expected false positives in 100 tests at  $\alpha=0.05$ ?
- How many expected false positives in 100 tests at FWER=0.05?
- Bonferroni
  - Instead of selecting 0.05 as your type I error rate, choose  $0.05 / \text{the number of tests}$
  - Makes no assumptions about correlation structure in the multiple tests
  - This is the most conservative way to deal with FWER, and is very hard to meet the significance threshold.



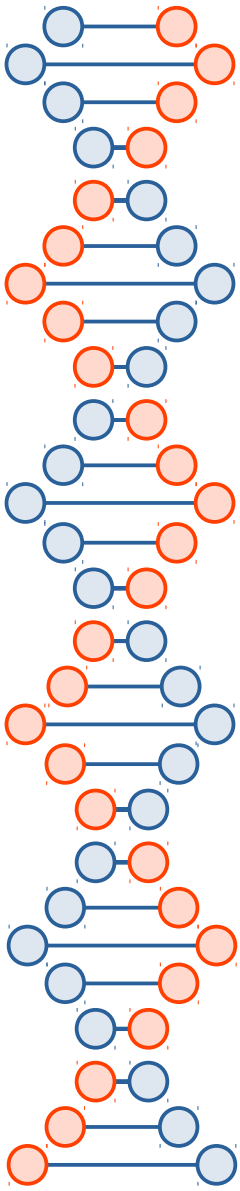
# Dealing with multiple comparisons

## Control for False Discovery Rate(FDR)

- Type I error rate:
  - In all the tests I've run, how many false positives were there?
- FDR:
  - Among the tests that were deemed significant, how many were incorrect.
- Benjamini-Hochberg Method
  - The least conservative technique and the most widely used in omics studies
  - Define  $q = V/R$  where
    - $V$  = # false positives,
    - $R$  = # discoveries (rejected  $H_0$ s).
  - Order the P-values in ascending order, smallest first:  
 $P(1) < \dots < P(N)$ .
  - For  $\alpha$ , find the largest  $k$  such that  $P(k) \leq (k/N)\alpha$ .
  - Then reject all  $H_0(i)$  for  $i = 1, \dots, k$ .

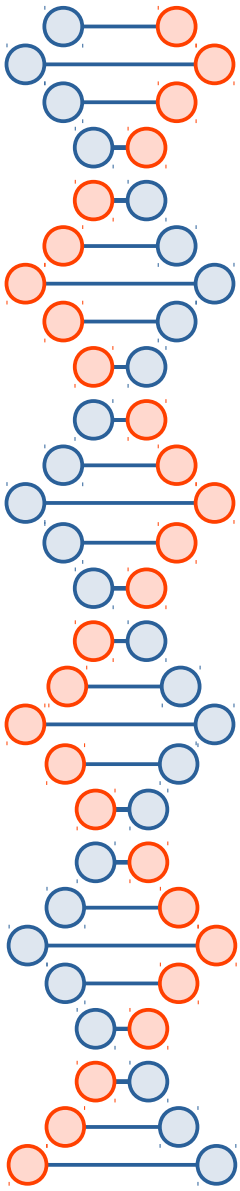






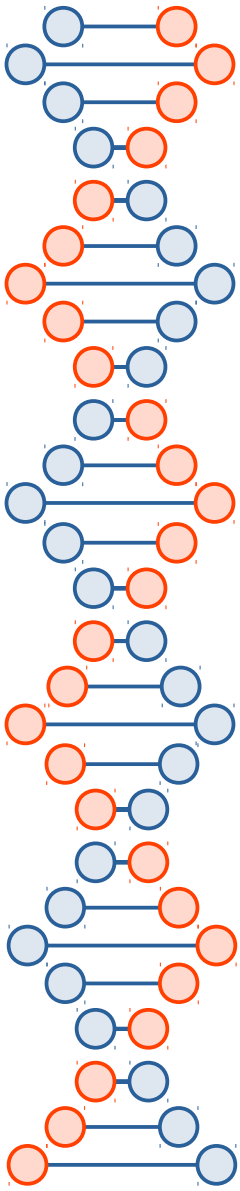
# Typical workflow in an Omics study

1. Plot global signal of all samples
2. Look for differential signals or associations across conditions
3. Make sure your method deals with multiple comparisons



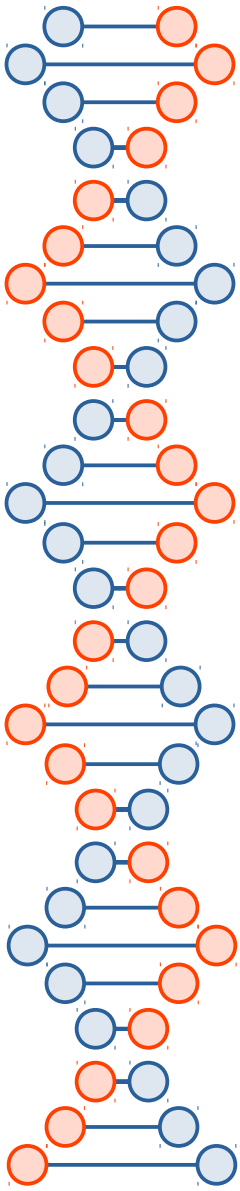
# Typical workflow in an Omics study

1. Plot global signal of all samples
  - Find outliers or potential problems with the way the data was prepared/processed
  - See if there are batch effects that needs to be addressed
  - See if there are differences between experimental conditions



# Typical workflow in an Omics study

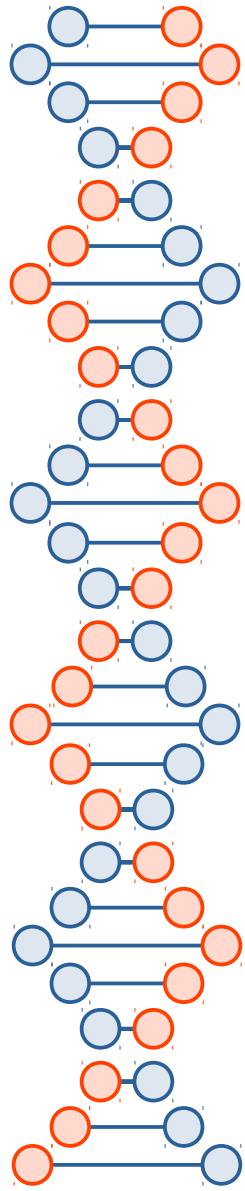
2. Look for differential signals or associations across conditions
  - Use an R package tailored to your protocol
    - `deseq2`
    - `LIMMA`
    - `methyKit`
    - more
  - Use specialized software developed for your protocol
    - `Homer`
    - `RMats`
    - `Bismark`
    - more
  - Write your own



# Typical workflow in an Omics study

## 3. Make sure your method deals with multiple comparisons

- Most of the well known bioinformatics software packages already take care of this
- At the very least, adjust for FDR



# Common filetypes for Omics

- What is a fastq file?
- What is a fasta file?
- What do you get when you map a fastq to a fasta?
- What is a bam file?
- What is a gtf file?
- What is a gene count file?