

CMSC423:

Bioinformatic databases, algorithms and tools

Héctor Corrada Bravo

Dept. of Computer Science

Center for Bioinformatics and Computational Biology

University of Maryland

University of Maryland, Fall 2014

Advances in Biology and Medicine needed, need, and will continue to need computational and statistical thinking (and their tools)

Héctor Corrada Bravo

Dept. of Computer Science

Center for Bioinformatics and Computational Biology

University of Maryland

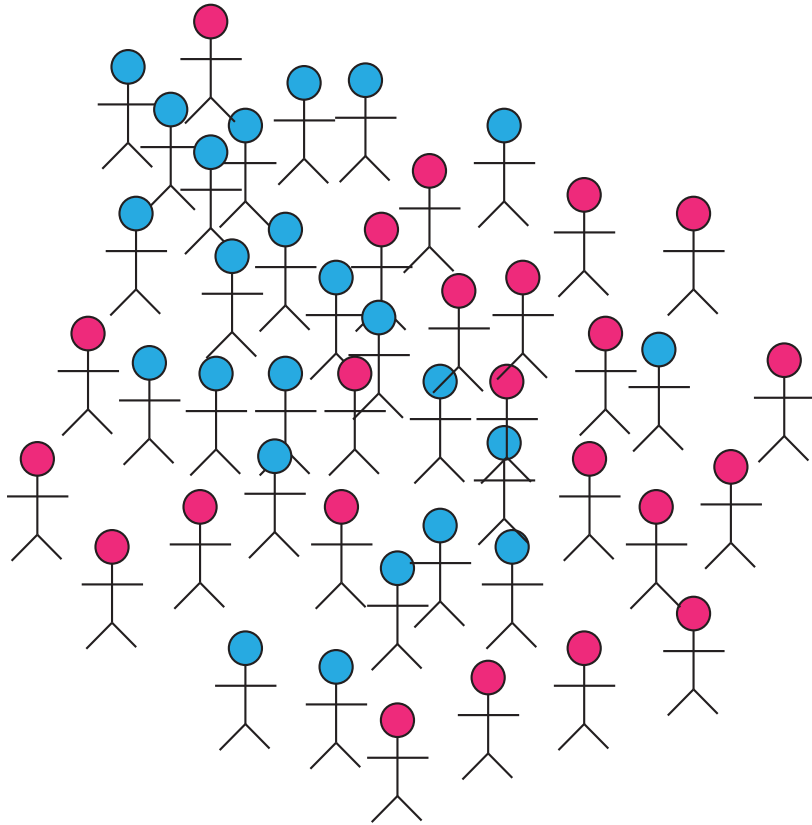
Why are my children
such pigs?



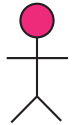
What is Genomics?

- Each cell contains a complete copy of an organism's **genome**, or blueprint for all cellular structures and activities.
- The genome is distributed along **chromosomes**, which are made of compressed and entwined **DNA**.
- Cells are of many different types (e.g. blood, skin, nerve cells), but all can be traced back to a single cell, the fertilized egg.

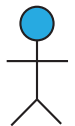
What is Genomics?



- Study the **molecular** basis of *variation* in development and disease
- Using **high-throughput** experimental methods
 - algorithms
 - ML
 - data management
 - modeling



cancer



healthy

Measurement

- For a small enough piece, we can measure the sequence of bases, referred to as *sequencing*
- Human Genome Project



D. melanogaster, Science, 2000



H. sapiens, Nature, 2000
and Science, 2000



M. musculus, Nature, 2002

Genome

[illegible]

**Total amount of DNA in human genome:
3 * 10⁹ base pairs (bp)**

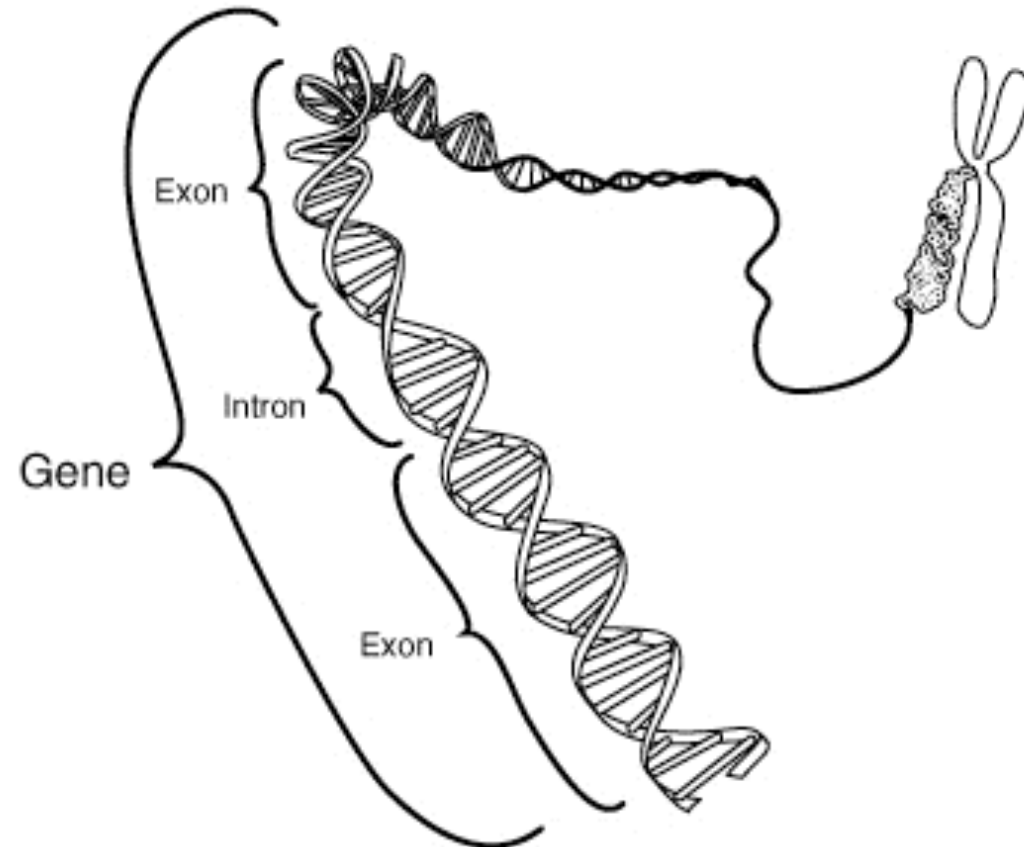
Why are these two different?



Differences explained by 1-10% difference in genome

Similarities explained by similar genes

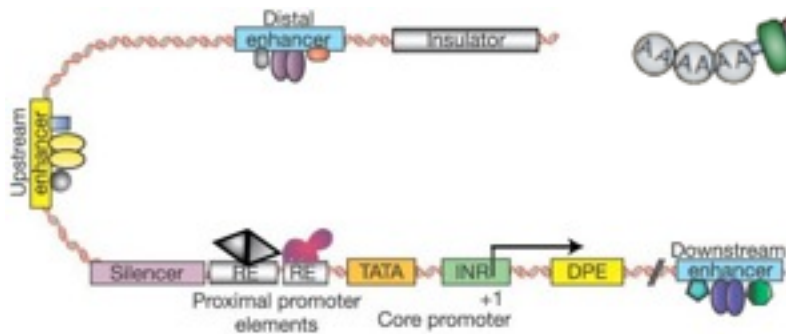
Genes



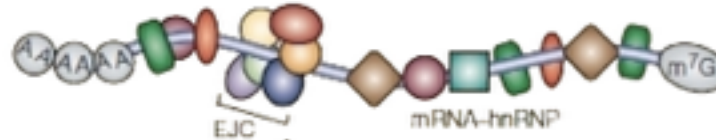
Computational Biology

Genes encode proteins which are transcribed into mRNA and translated into proteins.

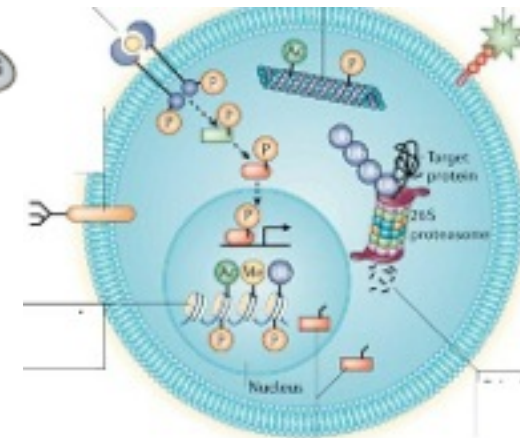
genomics



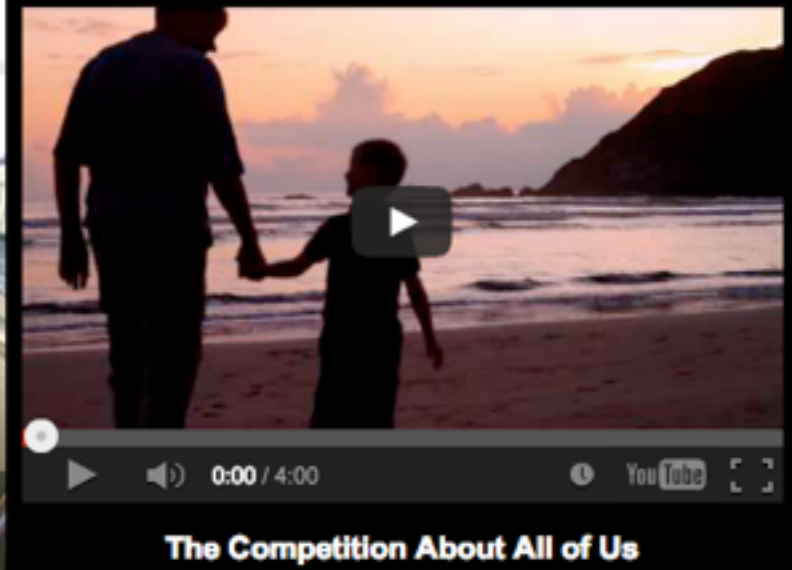
transcriptomics



proteomics



Major technological advances allow **unprecedented** data acquisition



build a **whole human genome** sequencing device and use it to sequence **100 human genomes** within **30 days or less**, with an accuracy of no more than one error in every 1,000,000 bases sequenced, with an accuracy rate of at least 98% of the genome, and at a recurring cost of **no more than \$1,000 (US) per genome**.

TOPICS

Space

Environment

Innovation

Weird science

COSMIC LOG

**\$10 million
Genomics X
Prize canceled:
'Outpaced by
innovation'**

CHINA

China gets set
to launch its
first moon
lander by
year's end

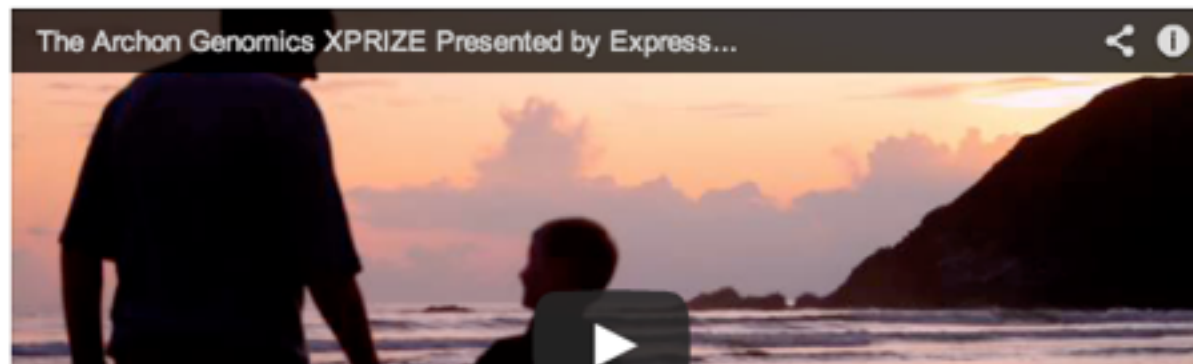
SPACE STATION

Japanese
astronaut to
command
space station in
March

\$10 million Genomics X Prize canceled: 'Outpaced by innovation'

Alan Boyle, Science Editor, NBC News

Aug. 23, 2013 at 11:11 PM ET



“genome sequencing technology is plummeting in cost and increasing in speed independent of our competition”

“companies can do this for less than \$5,000 per genome, in a few days or less — and are moving quickly towards the goals we set for the prize.”

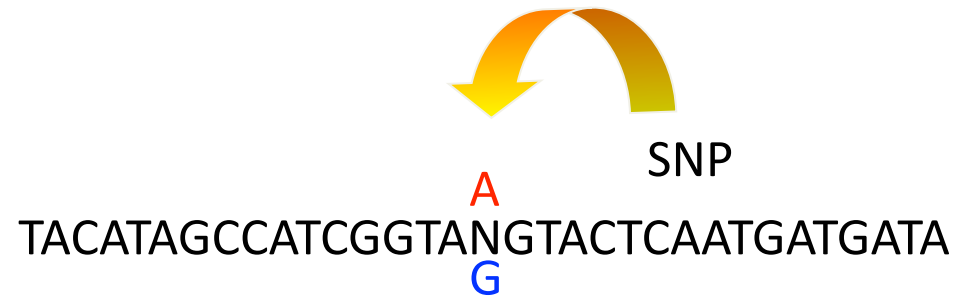
What makes them different?



Much human variation is due to difference in ~ 6 million base pairs (0.1 % of genome) referred to as SNPs

Single Nucleotide Polymorphism (SNP)

Genomic DNA:



From reads to evidence

```
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGCGCTGNNNNNNNNNNNNNNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A%NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>%NNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNNNNNNNNNNNNN
+
A=B7&7:>B@:A>?9:<;:>?4?%NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAGAAGCACAGCCAANNNNNANTNNTNNTNNN
+
BBCCCCCBB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNNGNTNAAAGNNNN
+
BCC7+<B=7BB5=ABA7B6BBB84BB7B%NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:947#0/1
CCAGGAGAAAGCCATGTTTCAGTTCGAGCGCNANANCGTGANNNN
+
BBB9@?7A7>AAB@>7B=?@.>87B?%NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:563#0/1
CCAGCCCCCTCCCCATCTCCACCCCTGTACCTNANCCCCGTGANNNN
+
BBABAABB;AAABA77@5AAA:??>%NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:1631#0/1
TGGGAACGCAGCCTACACTCTCCAGGCCTCCTNCCTCCGTNNNN
+
BBB@5@BBBBBBB@BBBBBAAABBB?;9BB@BA5&<B:%NNNNNN
@HWI-EAS146:5:1:2:1420#0/1
CTCAAACTCCTGACCTTGGTGATCCACCGCCTNGGCCCTCNNNN
+
BBBB:BBBBBABAAA?:(-AB@>AAA7AB7=A%NNNNNNNNNNNN
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGCGCTGNNNNNNNNNNNNNNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A%NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>%NNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNNNNNNNNNNNNN
+
A=B7&7:>B@:A>?9:<;:>?4?%NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAGAAGCACAGCCAANNNNNANTNNTNNTNNN
+
BBCCCCCBB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7CB7
```

From reads to evidence

I. Comparative

Sequence-wise, individuals of a species are nearly identical

Well curated, annotated “reference” genomes exist

```
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A\NNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGACGAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>\NNNNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNNNNCNNNNN
+
A=B767:>B@:A>79:<;>747\NNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAAGACAGCCAAANNNNNANTNNTCTNNNN
+
BBCCCCCBBBCB7CBC=7>+<=>=BCBCB\NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCTCINNNGNTNAAGNNNN
+
BCC?+<B=7BB5=ABA7B6BBB4BB7B\NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:947#0/1
CCCAGGAGAAAGCCATGTTCACTTCAGTCGAGCGCANNANANCGTANNNN
+
BBB9@77A7>AAB@>7B=?@.>B7B? \NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:563#0/1
CCAGCCCCCTCCCCATCTCCACCCCTGTACCTNANCCCCCTGANNNN
+
BBABAABB;AAABA77@5AAA:??\NNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:1631#0/1
TGGGAACGCAGCCTACACTCTTCCAGGCTCCTNCCTCCGTNNNN
+
BBB@6@BBB88888@BBBABAABBB?;9BB@BA5<B:\NNNNNNNN
@HWI-EAS146:5:1:2:1420#0/1
CTCAAACCTCTGACCTTTGGTGATCCACCCGCTNGGCCCTCNNNN
+
BBBB:BBBBBAAAA?:(=A@>AAA?AB?=\NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A\NNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGACGAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>\NNNNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNNNNCNNNNN
+
A=B767:>B@:A>79:<;>747\NNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAAGACAGCCAAANNNNNANTNNTCTNNNN
+
BBCCCCCBBBCB7CBC=7>+<=>=BCBCB\NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCTCINNNGNTNAAGNNNN
```



D. melanogaster, Science, 2000



H. sapiens, Nature, 2000
and Science, 2000



M. musculus, Nature, 2002



Idea: “Map” reads to their point of origin with respect to a reference, then study differences

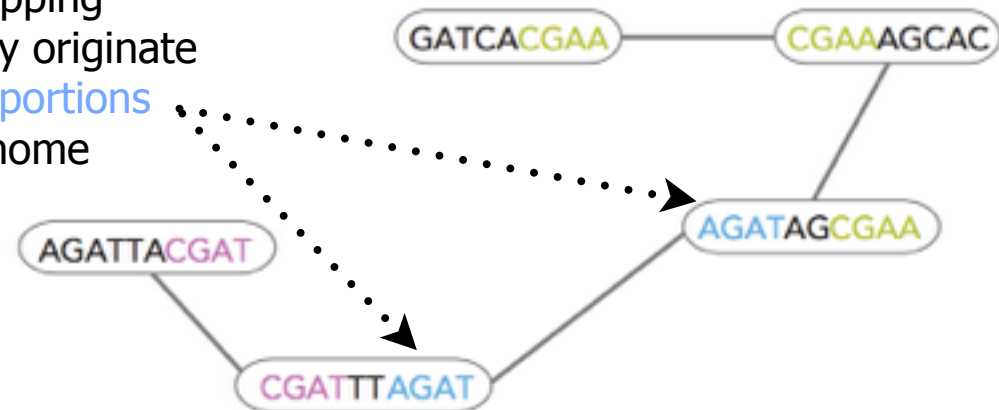
From reads to evidence

2. *de novo*

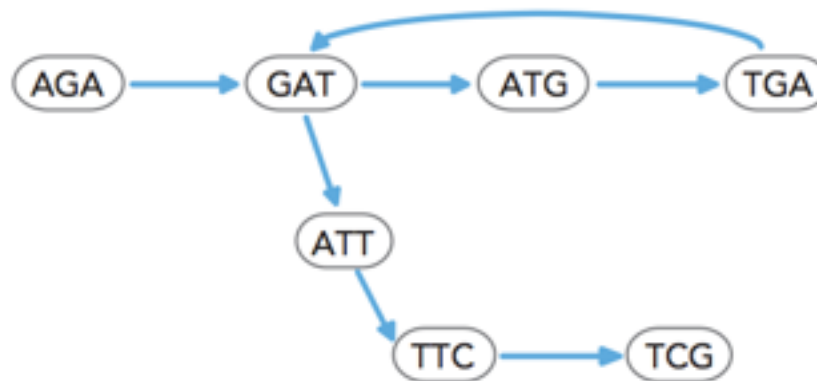
Assume nothing! - let reads tell us everything

Reads with overlapping sequence probably originate from **overlapping portions** of the subject genome

Encode overlap relationships as a graph



Source: De Novo Assembly Using Illumina Reads. Illumina. 2010



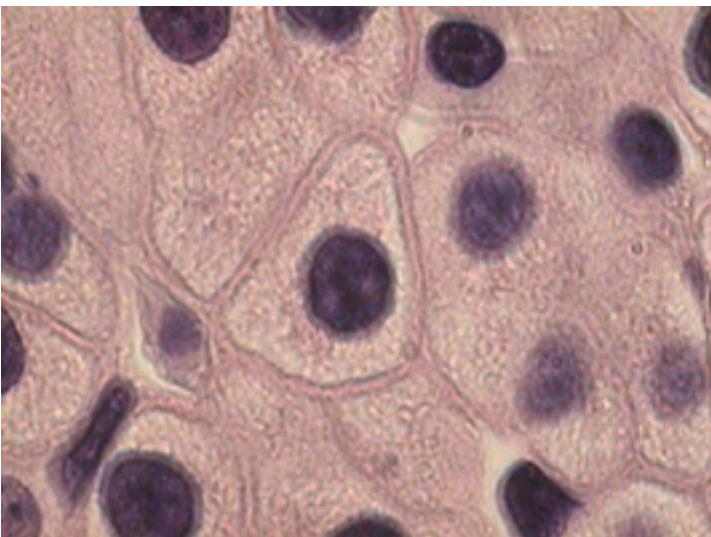
The full genome sequence is a "tour" of the graph

Source: De Novo Assembly Using Illumina Reads. Illumina. 2010

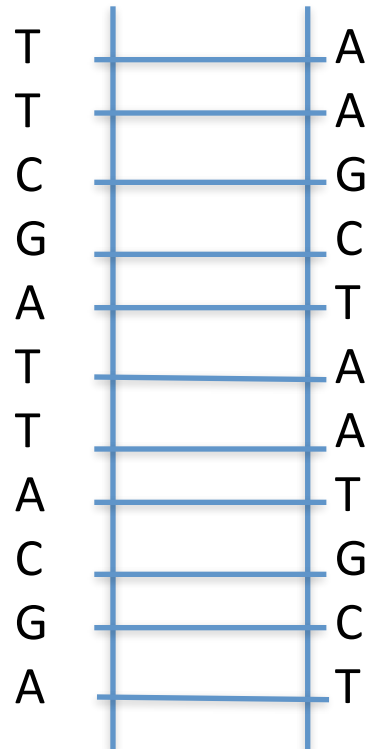
http://www.illumina.com/Documents/products/technotes/technote_denovo_assembly.pdf

```
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A\NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGACGAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>\NNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNNNNCNNNNN
+
A=B767:>B@;A>79:<;>747\NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAAGACAGCCAANNNNNANTNCTNNNNN
+
BBCCCCCBB7CBC=7>+<=BCBCB\NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCTCCHNNNGNTNAAGNNNN
+
BCC?+<B=7BB5=ABA7B6BBB4BB7B\NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:947#0/1
CCCAGGAGAAAGCCATGTTTCAGTTCGAGCGCHNANANCGTGANNNN
+
BBB9@77A7>AAB@>7B=7@.>B7B?NNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:563#0/1
CCAGCCCCCTCCCCATCTCCACCCTGTACCTNANCCCCCTGANNNN
+
BBABAABBB;AAABA77@5AAA:??\NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:2:1631#0/1
TGGGAACGCAGCCTACACTCTTCCAGGCTCCTNCCTCCGTNNNN
+
BBB@6@BBB88888@BBBABAABBB?;9BB@BA5<B:\NNNNNN
@HWI-EAS146:5:1:2:1420#0/1
CTCAAACCTCTGACCTTTGGTGATCCACCCGCTNGGCCCTCNNNN
+
BBBB:BBBBBAAAA?:(=AB@>AAA?AB?=\NNNNNNNNNNNN
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNNNNN
+
BBBB>A7B@;@BBBBBAA=BA=A\NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGACGAGGNNNNNNNNNNNGNNNNN
+
B9B@B<;BAA<@AB9=1>\NNNNNNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNNNNCNNNNN
+
A=B767:>B@;A>79:<;>747\NNNNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAAGACAGCCAANNNNNANTNCTNNNNN
+
BBCCCCCBB7CBC=7>+<=BCBCB\NNNNNNNNNNNNNNNNNN
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCTCCHNNNGNTNAAGNNNN
```

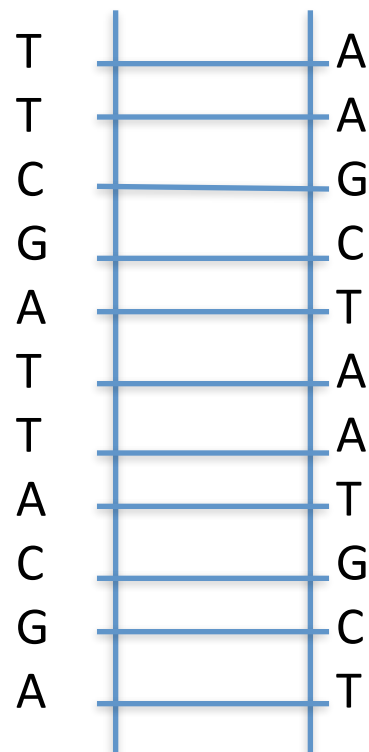
How many basepair differences?



Liver



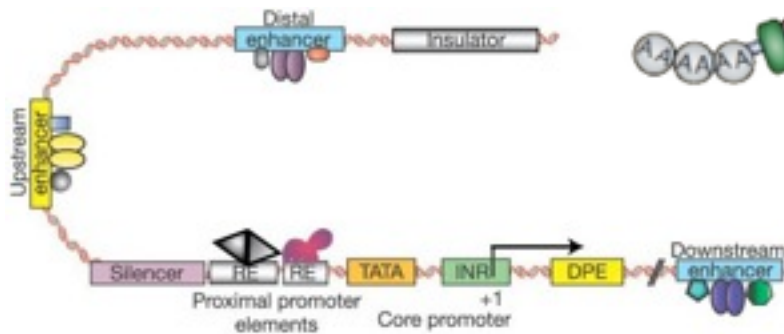
Brain



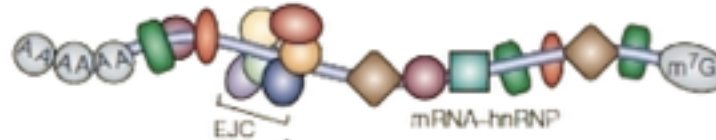
Computational Biology

Genes encode proteins which are transcribed into mRNA and translated into proteins.

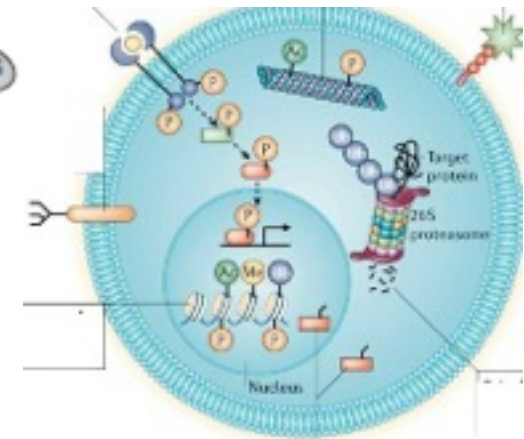
genomics



transcriptomics

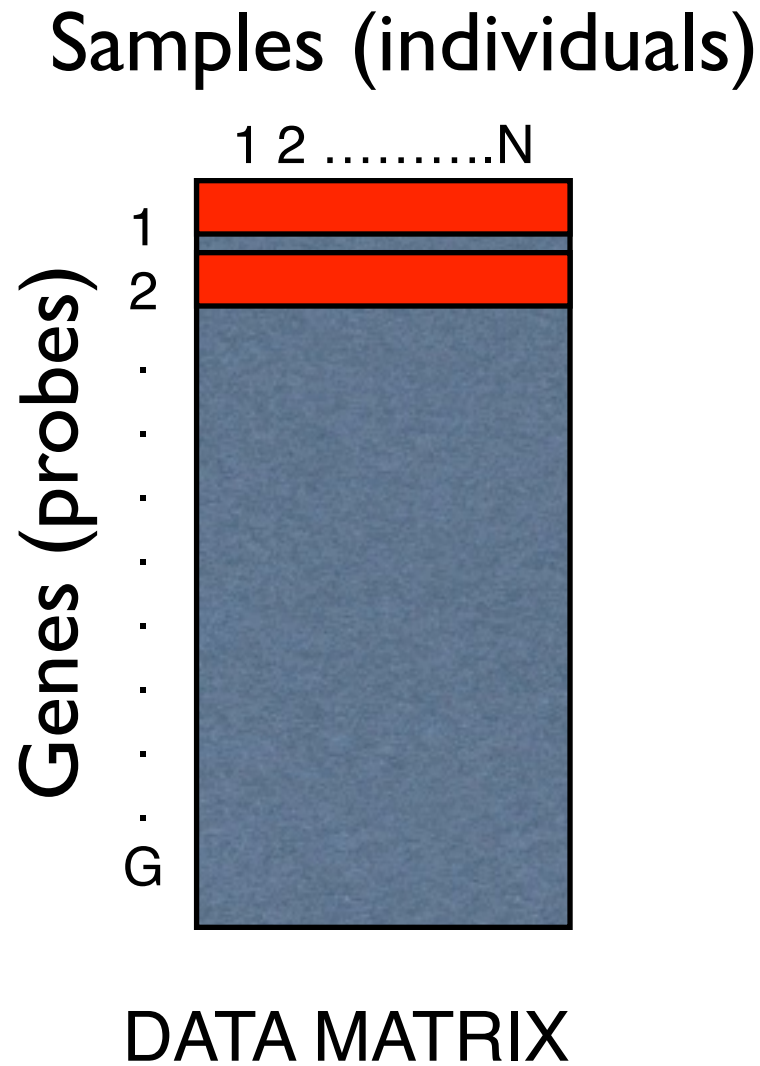


proteomics



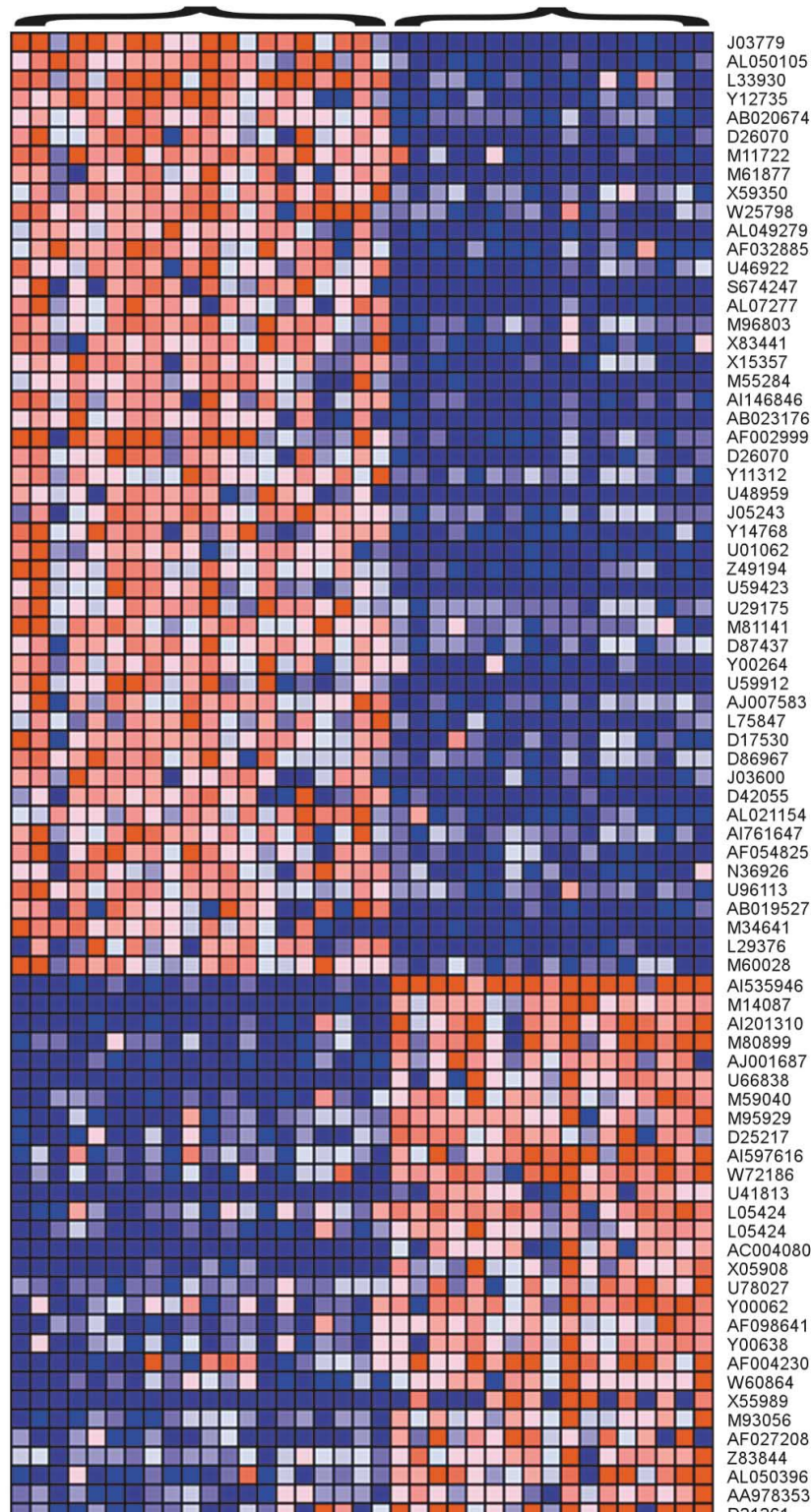
Major technological advances allow **unprecedented** data acquisition

Measurements



ALL

MLL



article

***MLL* translocations specify a distinct gene expression profile that distinguishes a unique leukemia**

Scott A. Armstrong¹⁻⁴, Jane E. Staunton⁵, Lewis B. Silverman^{1,3,4}, Rob Pieters⁶, Monique L. den Boer⁶, Mark D. Minden⁷, Stephen E. Sallan^{1,3,4}, Eric S. Lander⁵, Todd R. Golub^{1,3,4,5*} & Stanley J. Korsmeyer^{2,4,8*}

**These authors contributed equally to this work.*

Published online: 3 December 2001, DOI: 10.1038/ng765

Population Genomics

Clustering: Group samples (individuals) that show *similar* gene expression profiles

Classification: Discover gene expression profiles that distinguish two *populations*: e.g., cancer patients vs. healthy people

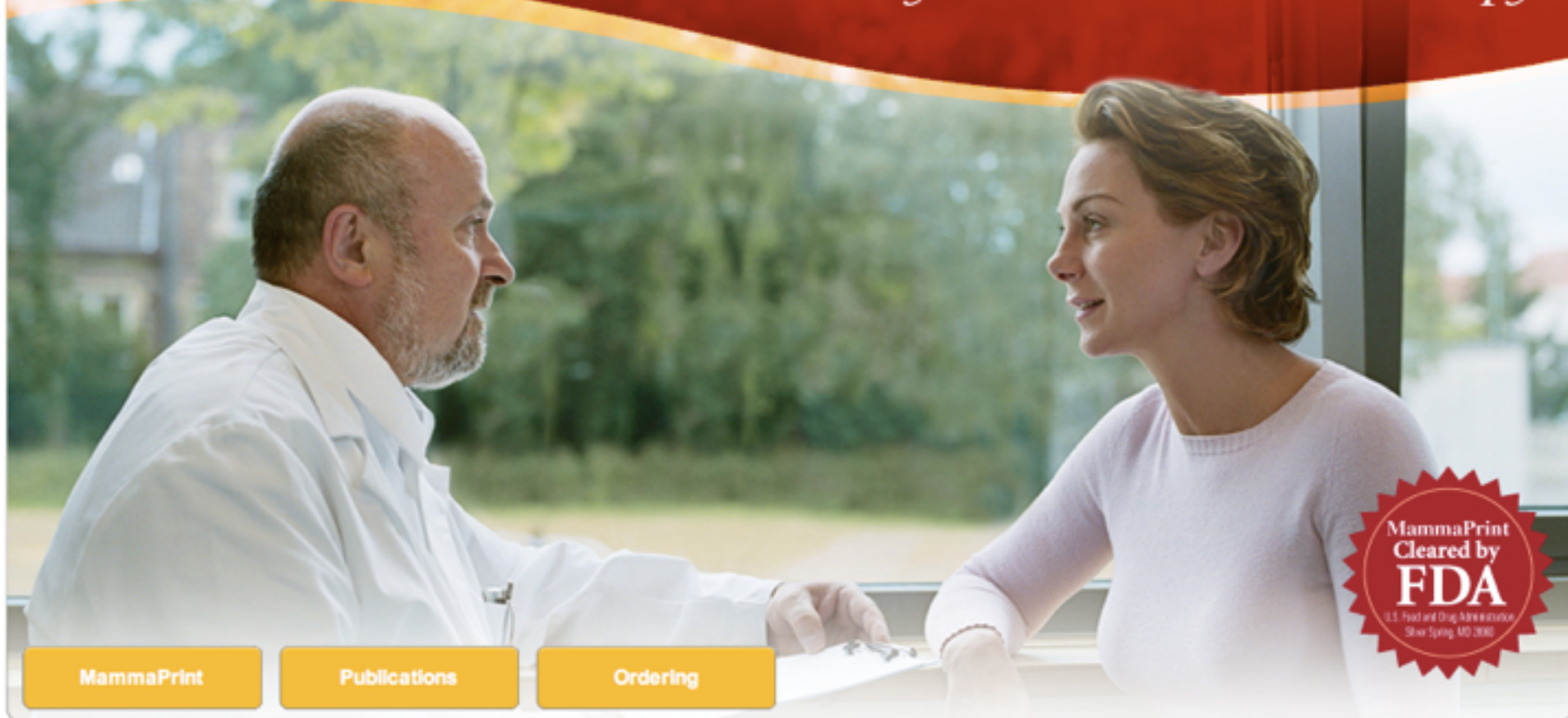
I.Networks: Discover groups of genes whose expression behaves differently in two populations

Why stats

If we want to infer things about gene expression in *populations*, we need to do some statistics

1. we want to see if some particular differences we see are due to *chance*
2. we want to make sure an experiment is setup so differences we see are those we care about
3. we want to have a sense of how general are inferences are (overfitting)

Learn how MammaPrint[®] can help
Personalize your breast cancer therapy.



MammaPrint

Publications

Ordering

What's New | April 2010: New Study Demonstrates Patients with High-Risk MammaPrint Profile Benefit from Chemotherapy

PERSONAL GENOMICS



23andMe genetics just got personal.

Go

[log in](#)

[claim codes](#)

[blog](#)

[help](#)

[your cart](#)

Get the latest on your DNA with \$399 and a tube of saliva



PERSONAL GENOMICS



- We need to produce reliable genome measurements, but on much bigger scale (**Algorithmics, Systems**)
- Multiple genome features, decide which are relevant and significant (**Information Retrieval, Data Management**)
- Population-based science, interpreted individually (**Machine Learning/ Statistics, Privacy**)

NHGRI strategic plan

- What does the NIH think genomics should be for the next 10 years?

PERSPECTIVE

doi:10.1038/nature09764

Charting a course for genomic medicine from base pairs to bedside

Eric D. Green¹, Mark S. Guyer¹ & National Human Genome Research Institute*

There has been much progress in genomics in the ten years since a draft sequence of the human genome was published. Opportunities for understanding health and disease are now unprecedented, as advances in genomics are harnessed to obtain robust foundational knowledge about the structure and function of the human genome and about the genetic contributions to human health and disease. Here we articulate a 2011 vision for the future of genomics research and describe the path towards an era of genomic medicine.

[Nature, Feb. 2011]

Where do we fit in?

- The major bottleneck in genome sequencing is no longer data generation—the computational challenges around data analysis, display and integration are now rate limiting. New approaches and methods are required to meet these challenges.
- **Data analysis**
 - Computational tools are quickly becoming inadequate for analysing the amount of genomic data that can now be generated, and this mismatch will worsen. Innovative approaches to analysis, involving close coupling with data production, are essential.
- **Data integration**
 - Genomics projects increasingly produce disparate data types (for example, molecular, phenotypic, environmental and clinical), so computational approaches must not only keep pace with the volume of genomic data, but also their complexity. New integrative methods for analysis and for building predictive models are needed.
- **Visualization**
 - In the past, visualizing genomic data involved indexing to the one-dimensional representation of a genome. New visualization tools will need to accommodate the multidimensional data from studies of molecular phenotypes in different cells and tissues, physiological states and developmental time. Such tools must also incorporate non-molecular data, such as phenotypes and environmental exposures. The new tools will need to accommodate the scale of the data to deliver information rapidly and efficiently.
- **Computational tools and infrastructure**
 - Generally applicable tools are needed in the form of robust, well-engineered software that meets the distinct needs of genomic and non-genomic scientists. Adequate computational infrastructure is also needed, including sufficient storage and processing capacity to accommodate and analyse large, complex data sets (including metadata) deposited in stable and accessible repositories, and to provide consolidated views of many data types, all within a framework that addresses privacy concerns. Ideally, multiple solutions should be developed^{[105](#)}.

Where do we fit in?

- Meeting the computational challenges for genomics requires scientists with expertise in **biology** as well as in informatics, **computer science**, **mathematics**, **statistics** and/or engineering.
- *A new generation of investigators who are proficient in two or more of these fields must be trained and supported.*

What else is the class about?

- Gives you an example of end-to-end use of what you've learned as CS as a practice
 - We discuss the design and analysis of algorithms (e.g., string algorithms, dynamic programming, iterative optimization methods)
 - We implement algorithms (python)
 - We analyze data (also in python)
- We also learn about biology, medicine and why government shutdowns are really awful

Administrative Details

Class webpage:

1. <http://www.cbcb.umd.edu/~hcorrada/CMSC423>

Everything you want to know is there.

1. Name
2. email (@umd.edu)
3. Department and degree
4. Are you registered?(Y/N)
5. Relevant CS background
6. Relevant stats background
7. Relevant biology background
8. What do you hope to get out of this class?
9. (a) Favorite, and (b) least favorite CS/stats term/name/word/phrase. Why?
10. (a) Favorite, and (b) least favorite biology term/name/word/phrase. Why?