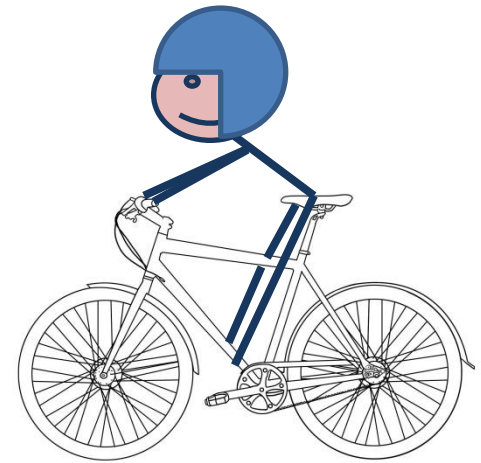# Genetic variation analysis: variant calling and annotations

Vincenza Colonna
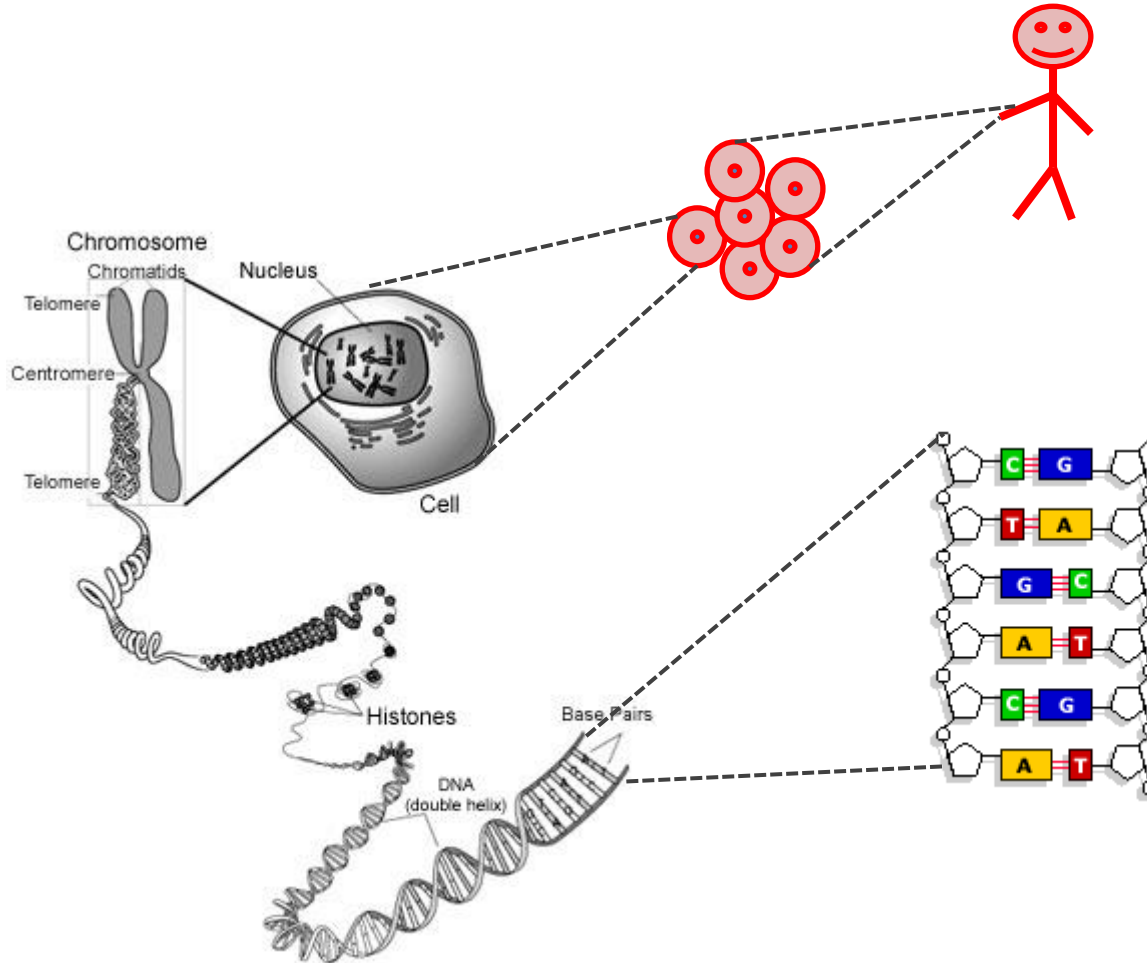
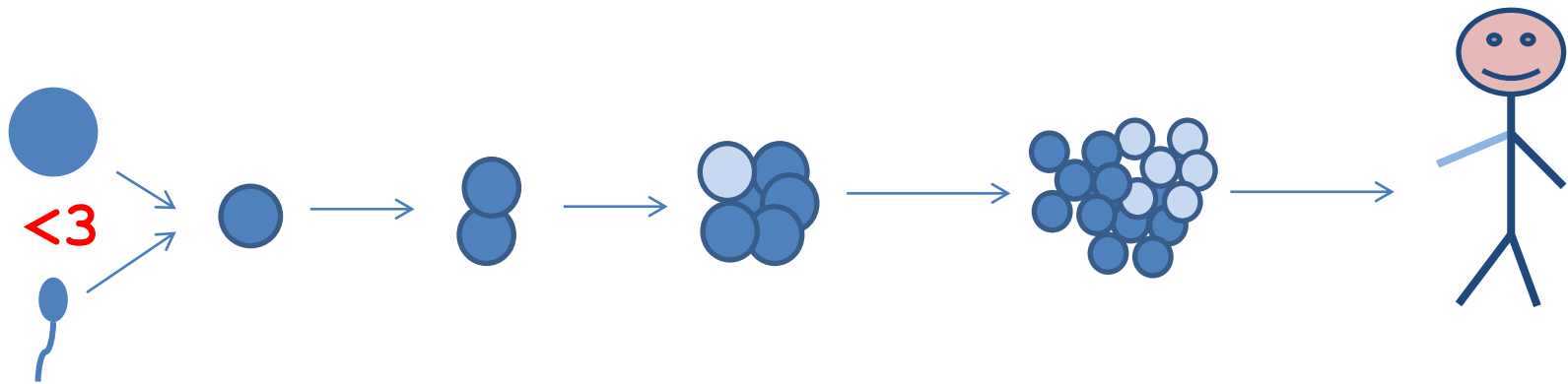**InterOmics Tutorial Day**

14 Novembre 2013

Area di Ricerca CNR, Via Castellino 111, Napoli

- Understanding the genomic variability in five minutes
- Few details on whole genome sequencing
- Variant detection – variant annotation
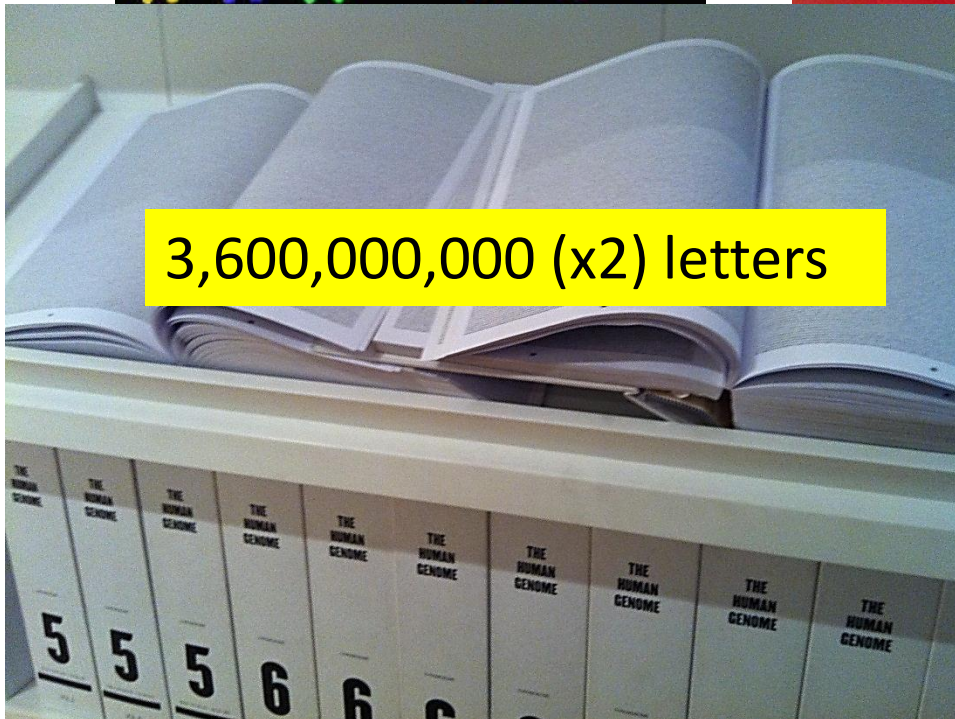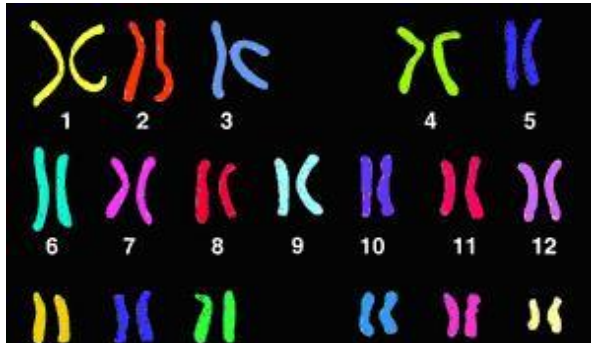- Practical session

# Where is the genome?

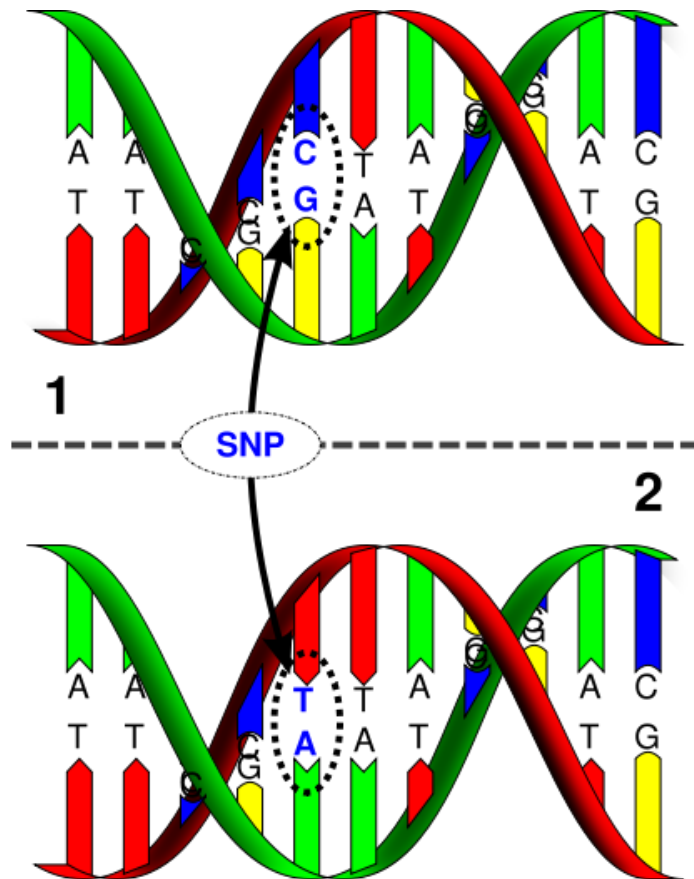# Does all the cells have the same genome in one organism?

**<3**

...well, mostly yes, but no...

# How big is the human genome?



3,600,000,000 (x2) letters

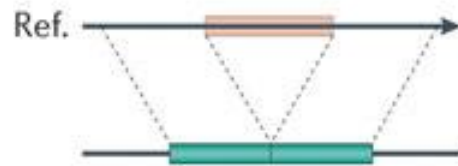# Is the DNA sequence identical among all genomes?



**SNPs**: Single letter changes

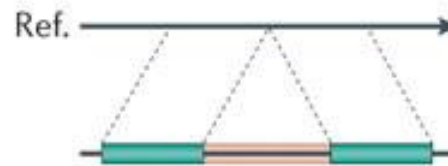**Indels**: Small insertions and deletions

**Structural variations**: Large changes in the structure and copy number of chromosomes or part of them

# Structural variants



Deletion

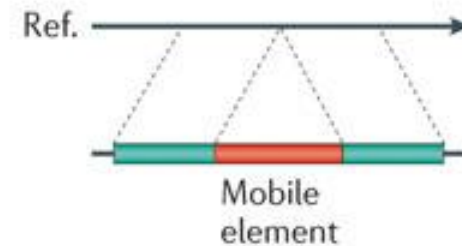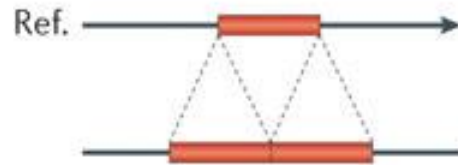Novel sequence insertion

Mobile-element insertion

Tandem duplication

Interspersed duplication

Inversion

Translocation

# Alleles or Variants

Allele 1

C  C  T  A

Allele 2

C  T  T  A

- Arise due to mutation

- Shuffled by recombination

Allele 3

C     T  A

- Diffused by migration

# Which are the consequence of DNA differences?

INTRA-SPECIES VARIABILITY

INTER-SPECIES VARIABILITY

DISEASES

- Understanding the genomic variability in five minutes

- **Few details on whole genome sequencing**

- Variant detection – variant annotation

- Practical session

# Chain Termination reaction

# Sequencing technology evolution



James Watson and Francis Crick

Walter Gilbert

Frederick Sanger

ILLUMINA
454
Solid

40 million clus
20 microns

| 1953 | 1977 | 1987 | 2007….. |
|---|---|---|---|
| | Maxam-Gilbert | Chain termination | High-throughput |
| | Sanger | | |

13

# "State of the art" technologies

Ion Torrent™ next-gen sequencing technology:

http://www.youtube.com/watch?v=MxkYa9XCvBQ

Pac Bioscience

http://www.youtube.com/watch?v=v8p4ph2MAvI

# How do we 'read' whole genomes?

1. DNA is extracted from donors and fragmented

*Many copies of the genome in fragments*

2. DNA sequence is determined for each fragment

**AATCTGTATG**
**TTCTGTC**
**ATTTCCTC**
**TTCAATC**

*MODERN "NEXT-GENERATION" SEQUENCING MACHINE*

3. Fragments are aligned against a reference sequence

**REFERENCE**
**IND1**

**AATCTGTATG**
**CCCGTAAAT     TATGCTTTT**

4. Overlapping fragments are merged into a 'consensus' sequence

**AATCTGTATG**
**CCCGTAAAT     TATGCTTTT**

*CCCGTAAATCTGTATGCTTTT*

# What do we sequence and for what?

✓ Variation → DNA-Seq **<3**

✓ Expression → RNA-seq

✓ Regulation → ChIP-seq

✓ Metagenomics → pooled DNA seq

✓ Non-model organisms…

# Why DNA-seq is so exciting?

**Cost per genome has decreased**



Adapted from NHGRI

**Sequencing Progress vs Compute and Storage**
Moore's and Kryder's Laws fall far behind

- Microprocessor (MIPS)
- Sequencing (kbases/day)
- Compact HDD storage capacity (MB)



Kahn (2011) *Science* **331**, 728-729

# Why DNA-seq is hard?

- Human genome is large:

  ~$3 \times 10^9$ nucleotides per haploid genome

- Sequence reads are short:

  35 – ~1,500 bp, with ~1% errors

- Understanding the genomic variability in five minutes
- Few details on whole genome sequencing
- Variant detection – variant annotation
- Practical session

# Variant discovery in a nutshell

**COVERAGE**

1. Read mapping

**REFERENCE**

2. Identification of variable sites

★ ★          ★          ★          **REFERENCE**

# Variant discovery in a nutshell

**1. Read mapping**

**2. Identification of variable sites**

**3. Individual genotype calling**

**4. Imputation**

COVERAGE

REFERENCE

REFERENCE

REFERENCE

HAPLOTYPES

# Variant calling algorithms

– Allele counting

– Probabilistic methods, e.g. Bayesian model

- quantify statistical uncertainty

- assign priors based on observed allele frequency of multiple samples

– Heuristic approach

- based on thresholds for read depth, base quality, variant allele frequency, statistical significance

# http://seqanswers.com/wiki/Software/list

few examples

1. http://samtools.sourceforge.net/mpileup.shtml
2. https://github.com/ekg/freebayes
3. http://www.broadinstitute.org/gatk/

# Discovering alleles using graphs (GATK HaplotypeCaller)



## A Read Layout

$R_1$: GACCTACA
$R_2$: ACCTACAA
$R_3$: CCTACAAG
$R_4$: CTACAAGT
A: TACAAGTT
B: ACAAGTTA
C: CAAGTTAG
X: TACAAGTC
Y: ACAAGTCC
Z: CAAGTCCG

## B Overlap Graph

## C de Bruijn Graph

Traverse the graph to enumerate the possible haplotypes. Each edge is weighted by the number of reads which gave evidence for that k-mer.

*Assembly of large genomes using second-generation sequencing. Schatz. Genome Research. 2010.*

13

*Courtesy of Erik Garrison*

# Haplotype detection (FreeBayes)

Reference

Reads

Detection window

Direct detection of haplotypes from reads resolves differentially-represented alleles (as the sequence is compared, not the alignment).

Allele detection is still alignment-based.

*Courtesy of Erik Garrison*

# Length-frequency spectrum



SNP (length=0) and indels

*Courtesy of Erik Garrison*

# Variant annotation!

Is there a functional consequence for a variant?



Chr 2 position 136608646: **T**

**Lactose tolerant**

Chr 2 position 136608646: **C**

**Lactose intolerant**

# How many variable sites are <u>expected</u> in the human genome?

- Mutation rate is ~$1.2 \times 10^{-8}$ bp$^{-1}$ generation$^{-1}$

- In every gamete ~ 30 bases mutate

- In a population of ~$7 \times 10^9$, almost every possible genetic variant will be present

# How many variable sites are observed in the human genome?

- Two human genomes typically differ at 3.5 millions of positions

- Thousand humans (two thousands genomes) typically have only 40 millions variable sites

→ Some sites can't change

# Consequences in coding
# (2% of the genome)

# Consequences in non coding (98% of the genome)

- Transcription factor binding sites

- Promoters

- Enhancers

- Chromatin modifications

- ….

# Consequences Prioritization

FunSeq

- http://funseq.gersteinlab.org/
- http://info.gersteinlab.org/FunSeq



E. Khurana, Y. Fu, V. Colonna, X. J. Mu, etal.. Integrative annotation of variants from 1092 humans: Application to cancer genomics. Science, 342(6154), 2013.

# Few examples of software for annotation

1. http://www.bioconductor.org/packages/2.13/bioc/html/VariantAnnotation.html

2. http://www.openbioinformatics.org/annovar/

3. http://vat.gersteinlab.org/www.openbioinformatics.org/annovar/

4. http://www.ensembl.org/info/docs/tools/vep/index.html  <3

# Sequence ontology terms

SIFT  http://sift.bii.a-star.edu.sg/

Polyphen  http://genetics.bwh.harvard.edu/pph2/

- Understanding the genomic variability in few minutes
- Few details on whole genome sequencing
- Variant detection – variant annotation
- Practical session

YU-HU!

# Thanks!

- ClaudiaR
- Mario
- Pasquale
- Valerio

*WIFIGB*

interomics

bioinformatica

*SSH*

[corso@bender.igb.cnr.it](mailto:corso@bender.igb.cnr.it)

bioinformatica

# Sequence file types

Table 1.2: Common file types
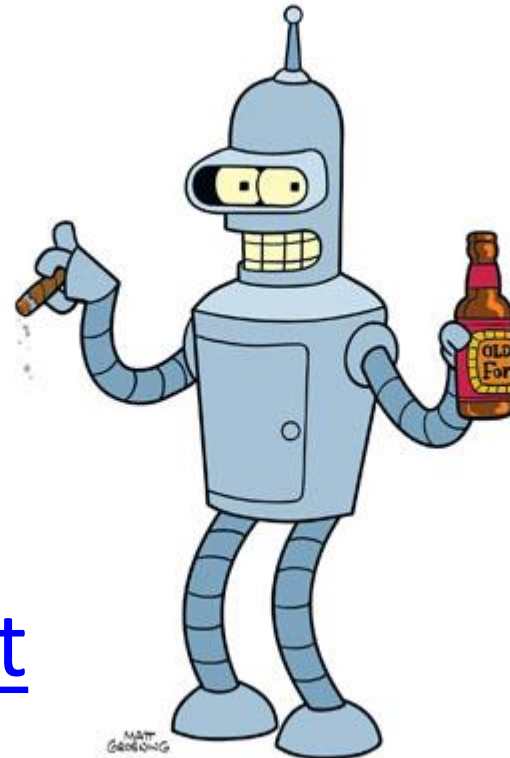
| File | Description |
| --- | --- |
| FASTQ | Unaligned sequences: identifier, sequence, and encoded quality score tuples |
| BAM | Aligned sequences: identifier, sequence, reference sequence name, strand position, cigar and additional tags |
| VCF | Called single nucleotide, indel, copy number, and structural variants, often compressed and indexed (with *Rsamtools* bgzip, indexTabix) |
| GFF, GTF | Gene annotations: reference sequence name, data source, feature type, start and end positions, strand, etc. |
| BED | Range-based annotation: reference sequence name, start, end coordinates. |
| WIG, bigWig | 'Continuous' single-nucleotide annotation. |
| 2bit | Compressed FASTA files with 'masks' |

# bgzip

- BAM files are compressed using a variant of GZIP (GNU ZIP), called BGZF (Blocked GNU Zip Format)

- BGZF is intended to improve on GZIP for random access.

# tabix

- Generic indexer for TAB-delimited genome position files

- fast retrieval of sequence features from a big tab-delimited file

# VCF, VCFtools, vcflib

- Poster
- http://www.1000genomes.org/wiki/Analysis/Variant%20Call%20Format/vcf-variant-call-format-version-40)
- Vcftools http://vcftools.sourceforge.net/
- Vcflib https://github.com/ekg/vcflib

# Variant effect predictor

- [http://www.ensembl.org/info/docs/tools/vep/index.html](http://www.ensembl.org/info/docs/tools/vep/index.html)
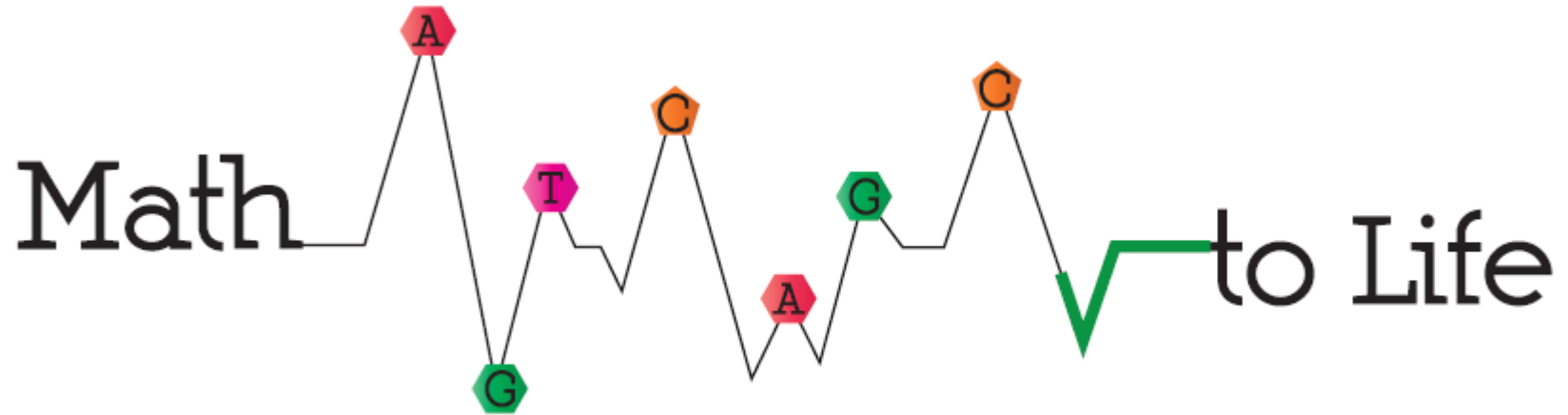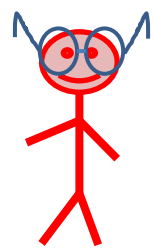

- ENCODE

- I have sequenced a number of individuals and I want to know allele frequencies in a subset of them

- I want to download 1000Genomes vcf file to use as comparison with the samples of Asians that I am studying

- I have discovered some variants in my samples of patients and I would like to know if there are functional consequences related to them

- I have discovered a variant in my samples and I would like to know if Neanderthal had it

# WORKSHOP ANNOUNCEMENT

Napoli May 2014

THAT'S INTERESTING

$$(x+a)^n = \sum_{k=0}^{n} \binom{n}{k} x^k a^{n-k}$$