

Tools and Algorithms in Bioinformatics
GCBA815, Fall 2013

Week-13: NextGen Sequence Analysis

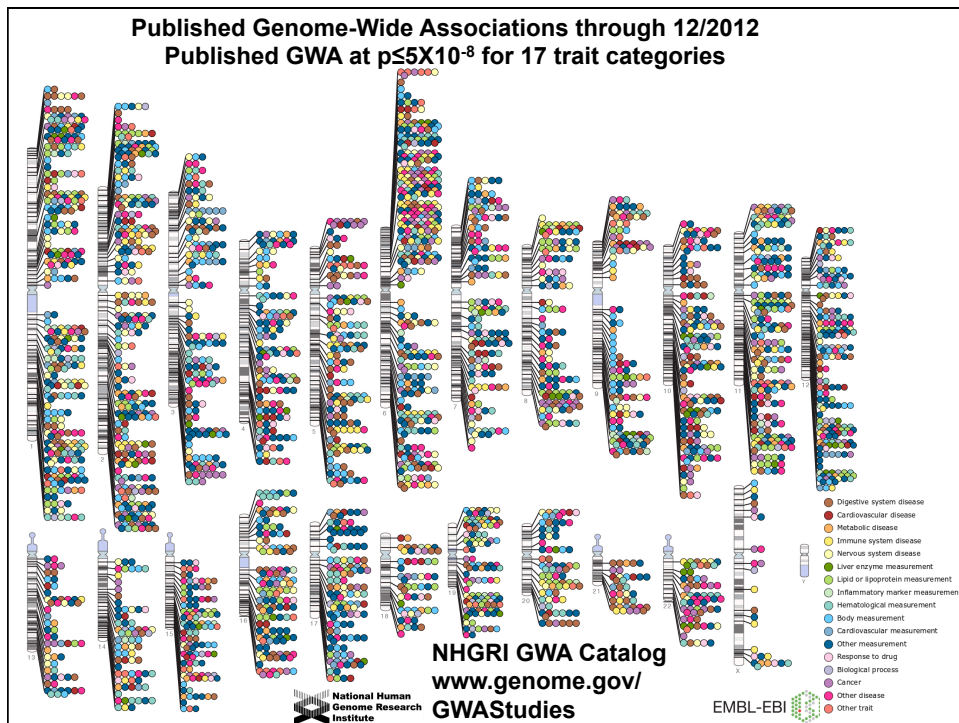
Demonstrators: Suleyman Vural, Li You, Sanjit Pandey

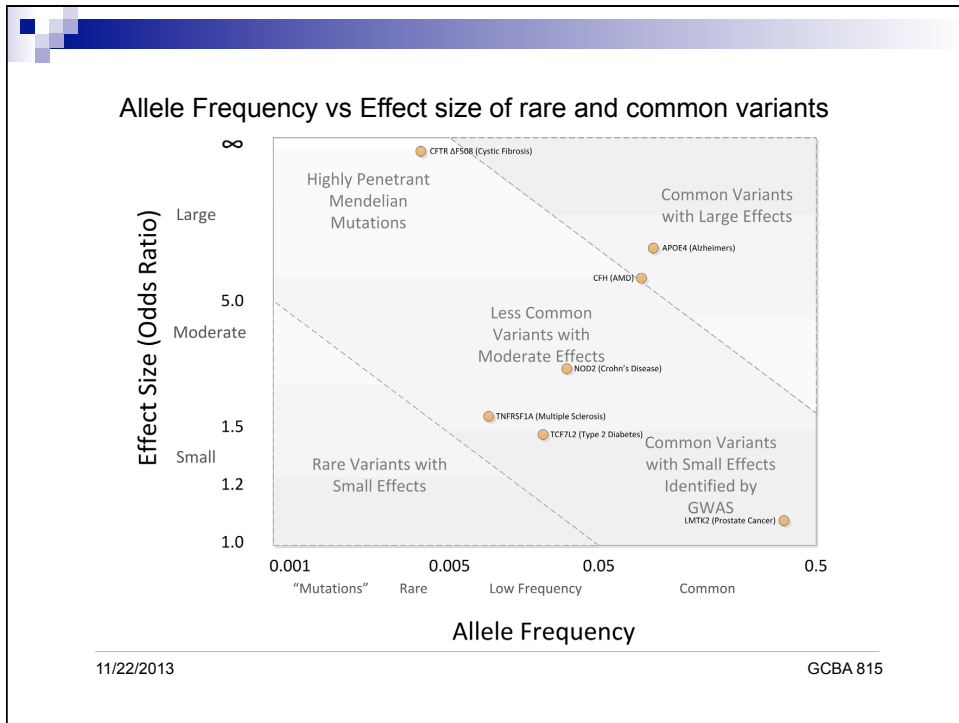
Babu Guda

Department of Genetics, Cell Biology and Anatomy
University of Nebraska Medical Center

11/22/2013

GCBA 815





IUPAC codes for nucleotides

| Code | Definition | Meaning |
|------|------------|-------------------|
| A | Adenine | A |
| C | Cytosine | C |
| G | Guanine | G |
| T | Thymine | T |
| R | AG | puRine |
| Y | CT | pYrimidine |
| K | GT | Keto |
| M | AC | aMino |
| S | GC | Strong |
| W | AT | Weak |
| B | CGT | Not A |
| D | AGT | Not C |
| H | ACT | Not G |
| V | ACG | Not T |
| N | AGCT | aNy |

11/22/2013

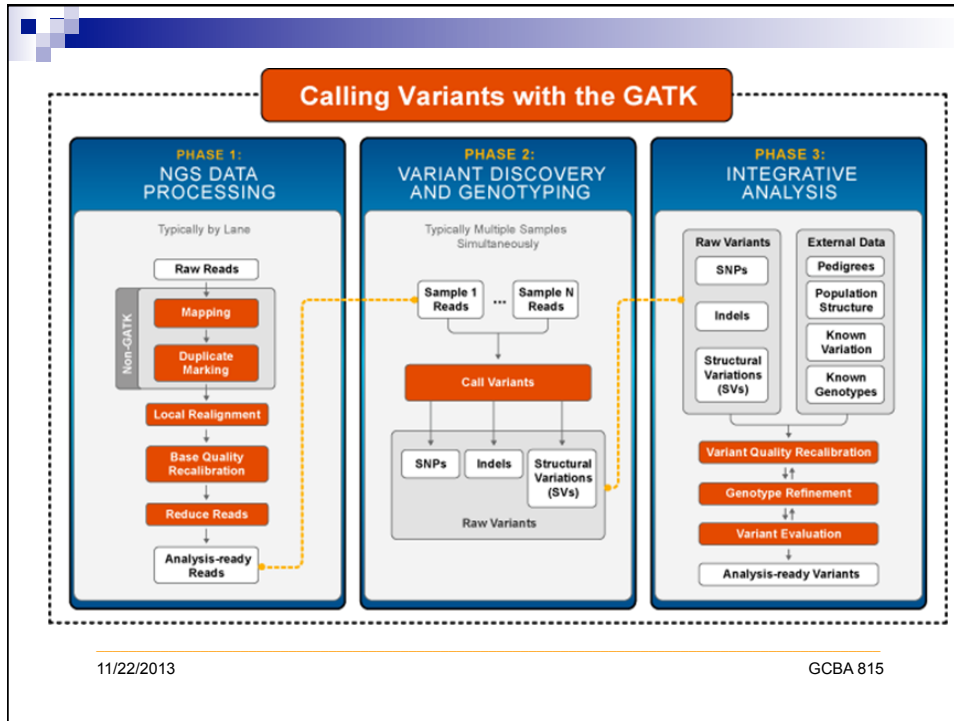
GCBA 815

NGS Applications in Biosciences

- Genome
 - Exome sequencing
 - Clinical sequencing, personalized medicine
 - Targeted genome sequencing (Ex: Ion torrent amplicons)
 - Whole genome sequencing
- Transcriptome
 - Whole transcriptome analysis
 - Small RNA analysis (siRNA, lncRNA, miRNA)
 - Gene expression profiling for selected target genes
- Metagenome
 - Sequencing together the genomes of a mixture of species
 - Example: Human gut microbiota or environmental samples
- Epigenome
 - Chromatin Immunoprecipitation Sequencing (ChIP-Seq)
 - Methylation and chromatin remodeling studies

11/22/2013

GCBA 815



- ### Different file types in NGS analysis
- Fastq file – generated by the sequencer, contains NGS reads
 - SAM file – Sequence Alignment/Map (generated by aligning the NGS reads with the reference genome)
 - BAM file – Binary version of the SAM file (SAMtools are used to manipulate SAM/BAM files)
 - GFF file – General Feature Format used to hold genome annotation (chromosome, strand, frame, exon, CDS, etc.)
 - GTF file – Gene Transfer Format (Also contains all the info as in GFF and in addition contains gene annotation information)
 - VCF file – Variant Call Format (used to store variant data such as SNPs, InDels, short structural rearrangements)
- 11/22/2013 GCBA 815

FASTQ format:

FASTQ is based on the popular FASTA format for sequences

FASTA format

>sequence_ID; header in one line
AGTTGTAGTCCGTGATAGTCGGATCGG

FASTQ format provides additional information that includes the quality score

```
@20FUKAAXX100202.1:64:10634:114560/1
TTGTATTTTAGTAGAGACGGAGTTTCGCCATGTTGGTCAGGCTGGCCTCGAATTCCTGACCTCAAGTGATCCGCCCGCCTCGGCCTCCCAACGTTTTGG
+
?=@7=>B==;BB?<B?=8539<676>8>=BB<<B=08:9@5:A@@?@9:BAAA<?:8:@AC@BBBBBA?<9-@B@:CAA77<:BEB<BB@07?@=<?84
```

ASCII code for Quality score (Phred score, ranges from 0-50)

ASCII code for Quality score (in the increasing order: ! is the worst and ~ is the best

```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNPOQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
```

11/22/2013

GCBA 815

Interpretation of Quality Score (Phred score)

Phred score (Q) vs Error probability (P)

$$Q = -10 \log_{10} P$$

Phred quality scores are logarithmically linked to error probabilities

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---------------------|------------------------------------|--------------------|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10000 | 99.99% |
| 50 | 1 in 100000 | 99.999% |

11/22/2013

GCBA 815

