

* EXONS -

- any part of a gene that will encode a part of the final mature RNA produced by that gene after introns have been removed by RNA splicing.
- refers to both the DNA seq. within a gene and to the corresponding seq. in RNA transcripts.
- How many exons does a gene have? - 8.8 exons.
- On an avg. there are 8.8 exons and 7.8 introns per gene.
- About 80% of the exons on each chromosome are < 200

* TOOLS -

- ① Genomescan
- ② Geneid
- ③ OrailEXP

* ORF [Open Reading Frame] -

- ① All ORF can be exons.
- ② All exons cannot be ORF.

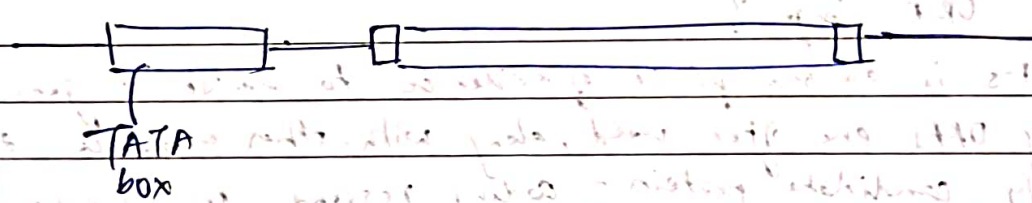
- portion of a DNA molecule that, when translated into a.a., contains no stop codons.
- A long ORF is likely part of a RF that has the ability to be translated.
- An ORF is a continuous stretch of codons that may begin with a start codon (usually AUG) and ends at a stop codon (usually UAG, UAA or UGA).
- Why are ORF imp?
 - ↳ ORFs is as one piece of evidence to assist in gene prediction.
 - ↳ Long ORFs are often used, along with other evidence, to initially identify candidate protein-coding regions or functional RNA-coding regions in a DNA sequence.

Tools - MCB1 - ORF FINDER.

ELECTIVES CONTINUED

- Identify ORFs -

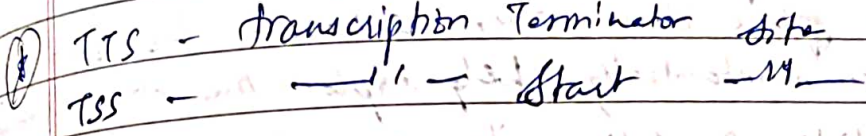
- simple 1st step in gene findings
- Translate genomic seq. in six frames.
- Identify stop codons in each frame.
- Regions without stop codons are called "open reading frame" or ORFs.
- Locate and tag all of the likely ORFs in a sequence.
- The longest ORF from a methionine codon is a good prediction of a protein encoding sequence.



* Promoters -

- The DNA opens up in the promoter region so that RNA polymerase can begin transcription.
- Each gene (or in bacteria, each grp of genes transcribed together) has its own promoter.
- A promoter contains DNA seq. that let RNA polymerase or its helper proteins attach to the DNA. i.e., A DNA seq. that the transcription apparatus recognizes and binds.

Usually found near the beginning of a gene, the promoter has a binding site for the enzyme used to make a mRNA molecule.



② Diagram - Transcription unit

* Tools (Eukaryotic + Prokaryotic Tools) -

- ① F-PROM
- ② TSSP
- ③ TSSW
- ④ TSSA
- ⑤ B-PROM

* Splice sites -

- A genetic alteration in the DNA seq. that occurs at the boundary of an exon and an intron (splice site).
- Splice sites are the seqs. immediately surrounding the exon-intron boundaries.
- Found - These sites are found at the 5' and 3' ends of introns. Most commonly, the RNA seq. that is removed begins with the dinucleotide GU at its 5' end, and ends with AG at its 3' end.
- i.e., the GU-AG rule (originally called the G⁺A⁻ rule in terms of DNA sequence).

- Why splice sites are imp?
 - o Mutations in these seqs may lead to retention of large segments of intronic DNA by the mRNA, or to entire exons being spliced out of the mRNA. These changes could result in prodⁿ of a non-functional protein.
 - o These donor sites, or recognition sites, are essential in the processing of mRNA.

Diagram Normal pre-mRNA

- Why splice sites prediction are imp?
 - o Prediction of splice sites where accurate localization of splice sites or substantially help explore the str. of genes.
 - o Accurate prediction of splice sites can setup the boundaries of exons which is critical in alternative splicing prediction.

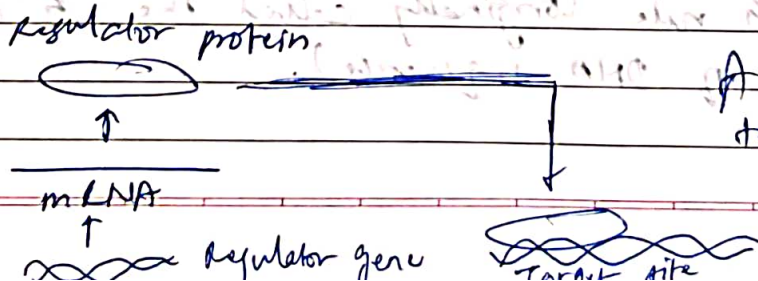
- ~~*~~ Tools (Eukaryotic + Prokaryotic Tools) (Abi based on various Approach / Algorithm) —

- ① Human Splicing Finder
- ② GenesSplicer
- ③ FGENES
- ④ Fgenesh-M
- ⑤ FGENESH-GC

* Regulatory Regions —

- A segment of a DNA molecule which is capable of increasing or decreasing the expression of specific genes within an organism.
- An enhancer activates the nearest promoter to it.
- A UAS [upstream activator seq.] in yeast behaves like an enhancer but works only upstream of the promoter.
- Form complexes of activators that interact directly or with the promoter.

Regulatory model



- 11
- ① the location of ORFs
 - ② Predict the structures of introns as well as exons of the genes of interest are of eukaryotic region.

PAGE No
DATE

What is the role of a regulatory region?

- Regulatory seq. controls when expression occurs for the multiple protein coding regions (red)
- Promoter, operator and enhancer regions (yellow) regulate the transcription of the gene into an mRNA.
- The mRNA untranslated regions (blue) regulate translation into the final protein products.

- Tools (Euk + Prokaryotic Tools) (also based on various approach / algorithm) -

- ① TRANSFAC ② QSAT Fungi / prokaryotes / bacteria
- ③ CRÈME ④ PSA Tools

* Prediction for various signals in genome is important?

- Why gene prediction?

With the rapid accumulation of genome seq info, there is a pressing need to use computational approaches to accurately predict gene structure.

* Computational gene prediction is a pre-requisite for detailed functional annotation of genes and genomes.

The process includes detection of -

- * The Ultimate goal -

- o To describe all the genes computationally with near 100% accuracy.
- o The ability to accurately predict genes can significantly reduce the amount of experimental verification work required.

- Disadvantage -

- o For eukaryotes, many problems in computational gene prediction are still largely unsolved.
- o This is because coding regions normally do not have conserved motifs.
- o The elements are diverse and not clearly defined.
- o Detecting coding potential of a genomic region has to rely on subtle features associated with genes that may be very difficult to detect.

* Normally elements are short (6-8 nucleotide) and found in any seq. by ~~random~~ ^{random} chance, thus ^{high} rate of false (+ve) results because of which sensitivity drops and specificity is hampered.

- solⁿ -

For preliminary identification of these elements → combine a multitude of features and use sophisticated algorithms that give either -

- ① ab initio - based predictions OR
- ② Predictions based on evolutionary info. (homology based) OR
- ③ Experimental data.

* Types of Approaches -

① Ab initio - based (prediction based on given seq. only) -
Relying on 2 major features associated with genes -

1st feature - is the existence of gene signals, which include start and stop codons, intron splice signals, transcription factor binding sites, ribosomal binding sites at polyadenylation (poly-A) sites, the triplet ^{codon} ~~codon~~ etc.

- Thus, unique features can be detected by applying probabilistic models such as Markov models or hidden Markov models, HMM to help distinguish coding from non-coding regions.

② Homology-based approaches

- Predictions based on significant matches of the query sequence with sequences of known genes.

eg. if a translated DNA seq is found to be similar to a known

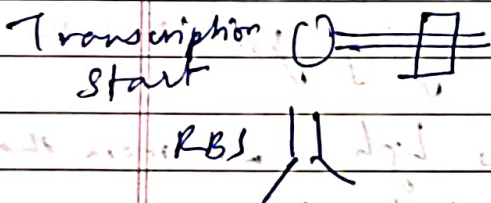
Therefore, there are also a no. of programs that actually combine prediction results from multiple individual programs to derive a consensus prediction. This type of algorithms can therefore be considered as consensus-based.

* Gene Prediction in Prokaryotes -

- Prokaryotes, which include bacteria and Archae, have relatively small genomes with sizes ranging from 0.5 to 10 Mbp ($1 \text{ Mbp} = 10^6 \text{ bp}$).
- The gene density in the genomes is high, with more than 90% of a genome sequence containing coding seq.
- There are very few repetitive sequences.
- Occasionally ATG and TTG are used as alternative start codons, but methionine is still the actual a.a. inserted at the first position.

- There may be multiple ATG, GTC or TAT codons in a frame.
- But presence of these codons at the beginning of the frame does not necessarily give a clear indication of the translation initiation site.
- Instead, to help identify this initiation codon, other features associated with translation are used.

- In many bacteria, it has a consensus motif of AGGAGG.
- Identification of the ribosome binding site can help locate the start codon.



- At the end of the protein coding region is a stop codon that causes translation to stop.
- Many prokaryotic genes are transcribed together as 1 operon.
- The end of the operon is characterized by a transcription termination signal \rightarrow ρ - independent terminator.

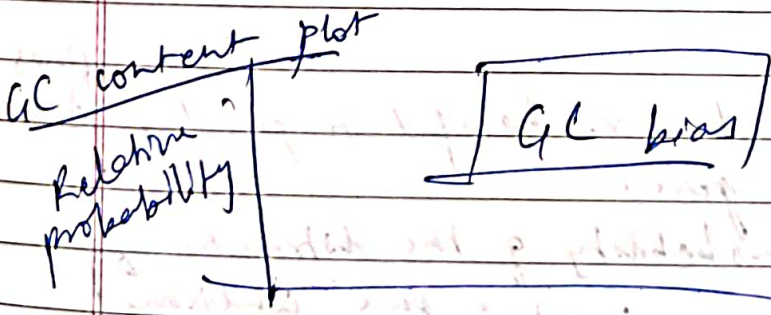
* Conventional Determination Method of ORF -

- Without use of specialized programs, prokaryotic gene identification can rely on manual determination of ORFs and major signals related to prokaryotic genes.
- Prokaryotic DNA is ~~not~~ first subjected to conceptual translation in all 6 possible frames, three frames forward and three frames reverse.
- Because a stop codon occurs in about every 20 codons by chance in a non-coding region, a frame longer than 30 codons without interruption by stop codons is suggestive of a ~~not~~ gene coding region.

HK

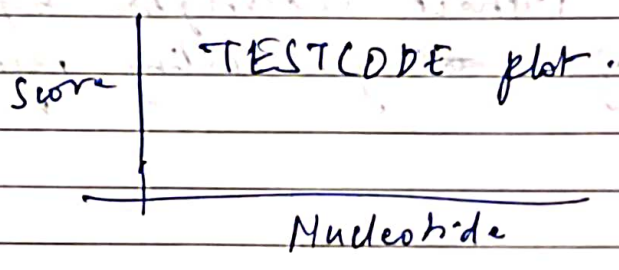
- In the early stages of development of gene prediction algorithms, genes were predicted by examining the non-randomness of nucleotide distribution.
- One method is based on the nucleotide composition of the 3rd position of a codon. Thus, in coding seq, it has been observed that this position has a preference to use G or C over A or T. (⊗) Koobke hypothesis.

Mk



- In addition to codon bias, there is a similar method called TESTCODE that exploits the fact that the 3rd nucleotide in a coding region tend to repeat themselves.
- By plotting the repeating patterns of the nucleotides at this position, coding and non coding regions can be differentiated.

NK



- Thus, statistical methods, which are based on empirical rules, examine the statistics of a single nucleotide (either G or C).
- They identify only typical genes and tend to miss atypical genes in which the rule of codon bias is not strictly followed.

NK

* Gene Prediction Using Markov Models and Hidden Markov Models -

- MM and HMM can be v. helpful in providing ^{finer} ~~more~~ statistical description of gene.
- MM - describes the probability of the distribution of nucleotides in a DNA seq, in which the conditional probability of a particular seq. position depends on k previous positions.

- In this case, k is the order of n MM.

NK

- A second-order model looks at the preceding 2 bases to determine which base follows, which is more characteristic of codons in a coding seq.

- The use of MM in gene finding exploits the fact that oligonucleotide distribution in the coding regions are diff. from those for the non-coding regions.

- These can be represented with various orders of MM.

- Since a fixed-order Markov chain describes the probability of a particular nucleotide that depends on previous k nucleotides, the longer the oligomer unit, the more non-randomness can be described for the coding region.

NK

- Once the parameters of the model are established, it can be used to compare the non-random distributions of trimers or hexamers in a new seq. to find the regions that are compatible with the statistical profiles in the learning set.

Gene prediction - Xin Xiong

- The frequency of ~~the~~ 6 unique nucleotides appearing together in a coding region is much higher than by random chance.
- ∴ a fifth-order Markov, which calculates the probability of hexamer bases, can detect nucleotide correlations found in coding regions more accurately.
- Problem of fifth-order Markov chain - is that if there are not enough hexamers which happens in short gene sequences, the method's efficacy may be limited.

NK

- Sometimes, genes tend to escape detection using the typical gene model. Thus, to make the algorithm capable of fully describing all genes in a genome, more than one MM is needed and therefore HMM prediction algorithms are implemented

* Performance Evaluation -

Accuracy of a prediction program - Sensitivity & Specificity
- calculated on 4 features accurately which are -

- ① TP - correctly predicted feature
- ② FP - incorrectly predicted feature
- ③ FN - missed feature
- ④ TN - correctly predicted absence of a feature.

② specificity (Sp) = $TP / (TP + FP)$ (Proportion of true signals among all signals that are predicted i.e., an ability to exclude incorrect predictions).

- A program is considered accurate if both sensitivity and specificity are simultaneously high and approach a value of 1.
- If the sensitivity is high but specificity is low, the program is said to have a tendency of over predict.

NOTE

- In the field of gene finding, a single parameter known as the correlation coefficient (CC) is often used, which is defined by

$$CC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TN + FN)(FP + FN)(TP + TN)}}$$

- CC provides an overall measure of accuracy, which ranges from -1 to +1, with
 - +1 = always correct prediction
 - 1 = incorrect prediction

* Gene Prediction in Eukaryotes -

NOTE

- The space b/w genes is often very large and rich in repetitive seqs and transposable elements.
- Eukaryotic genomes are characterized by a mosaic organization in which a gene is split into pieces (exons) by intervening non-coding seqs (introns).
- A eukaryotic gene is modified ~~to~~ in 3 diff. ways before becoming a mature mRNA for protein translations.
 - Capping at the 5' end of the transcript (methylation at the ~~last~~ initial residue of the RNA).

NR

- Computationally very demanding because of the presence of split gene structures, alternative splicing and very low gene densities.

100% - some conserved seq. features in eukaryotic genes that allow the computational prediction.

Eg - The splice junctions of introns and exons follow the GT-AG rule in which an intron at the 5' splice junction has a consensus motif of GTAACT, and at the 3' splice junction is a consensus motif of (Py)₁₂NCAG.

NR

- In addition, most of these genes have a high density of CG dinucleotides near the transcription ~~start~~^{start} site.
- And this region is referred to as a CpG island (p refers to the Phosphodiester bond connecting the 2 nucleotides), which helps to identify the transcription initiation site of a eukaryotic gene.

The poly-A signal can also help locate the final coding sequence.

Gene Prediction Programs -

- 3 categories of algorithms
- ① Ab initio based
 - ② Homology based
 - ③ Consensus based

- ① a. Prediction using Neural Network
- b. Using Discriminant Analysis
- c. Using HMM's.

①

a. Gene signals - includes signals such as gene start & stop sites and putative splice sites, recognizable consensus seqs such as poly-A sites.

b. Gene content - includes coding statistics, such as non-random nucleotide distribution, a.a. distribution, synonymous codon usage and hexamer frequencies

Neural Network / ANN for gene prediction -

Another feature that it resembles to nervous system is - its ability to "learn" and then make predictions after being trained.

The network is able to process info. and modify parameters of the weight functions & variables during the training stage. Once it is trained, it is able to make automatic predictions about the unknown.

// A NN is constructed with 3 layers

NK

- ↳ The gene str. info. is separated into several classes of features such as hexamer frequencies, splice sites and GC composition during training.
- ↳ The weight functions in the hidden layers are adjusted during this process to recognize the nucleotide patterns and their relationship with known structures.

NK

- Tool - GRAIL

① b.

~~that best separates~~

// and ∴ any identified variation could be responsible for problems.

Genome Analysis, Assembly & Annotation Studies

PAGE No

DATE

* Synteny -

(HK)

Use - provides a framework in which conservation of homologous genes and gene order is identified b/w genes of diff. species.
- The availability of human and mouse genomes paved the way for algorithm development in large scale called as synteny mapping, which eventually became an integral part of comparative genomics.

* What is synteny mapping?

(HK)

- Gene order conservation is in fact rarely observed among divergent species.
- \therefore comparison of syntenic relationships is normally carried out b/w relatively close lineages.
- However, if syntenic relationships for certain genes are indeed observed among divergent prokaryotes, they often provide imp. clues for functional relationships of the genes of interest.

eg - Genes involved in the same metabolic pathway tend to be clustered among phylogenetically diverse organisms. The preservation of the gene order is a result of the selective pressure to allow the genes to be coregulated and function as an operon.

* Reference genome seq. / Reference genome assembly -

(HK)

eg - The human reference genome doesn't represent the genetic seq. for any one individual, but is made up of a combination of several people's DNA. When we sequence a patient's genome and compare it to the reference genome, we assume that the reference represents the 'normal' sequence ...

Contd. →

Expensive technique (disadv).

How to search - (NK)

① * Gene expression profiling - (NK)

- Methods for analysis -

- ① Microarray #
- ② RT-PCR
- ③ qPCR
- ④ Northern blot
- ⑤ ELISA Test
- ⑥ Western Blot etc.

→ Advantage -

- ① Understand molecular mechanisms of complex disorders.
- ② Cancer studies
- ③ Metabolic disorder studies etc.

* Structural variants of DNA

- ① Relevance in molecular level processes
- ② Identification
- ③ Assembly of data from genome sequencing

- Aptamers and gene order - Any mol. bio book

- Mutation and genetic disease - refers to chromosomal studies
↳ case study approach. (Mendelian inheritance problem)

- Integrated genomic maps - proposed

for OMIM database
Genomic + Proteomics

or
only Genomics

or
only Proteomics

Om, all theory, methodology
O/P of that research.

* Structural Variants -

(NK)

- Microscopic & submicroscopic structural variant such as deletions, duplications, large copy no. variants as well as insertions, inversions and translocations.

* What is Variant ^{DNA} ~~data~~ -

(NK)

- Variant can be used to describe an alteration that may be pathogenic / of unknown significance.
- May be germline or somatic.

- Why are structural variants imp?

(NK)

- ↳ Regulation of gene expression, ethnic diversity and large-scale chromosome evolution - giving rise to the differences within populations and among

- How do you identify structural variations -

(NK)

- It is widely agreed that the method ^{can be classified} into 4 algorithms - Read-pair (RP), Split-read (SR), Read-Depth (RD) and Assembly

- What seq. tech. can be used to detect SV?

(NK)

- Strand-seq is the most suitable detection method for chromosomal inversions, a particularly challenging type of SV.

* DNA sequencing -

- ① Maxam-Gilbert
- ② Sanger - key reagent - ddNTP
- ③ Restriction endonucleases - Bacterial enzymes that cleave DNA at specific sites.

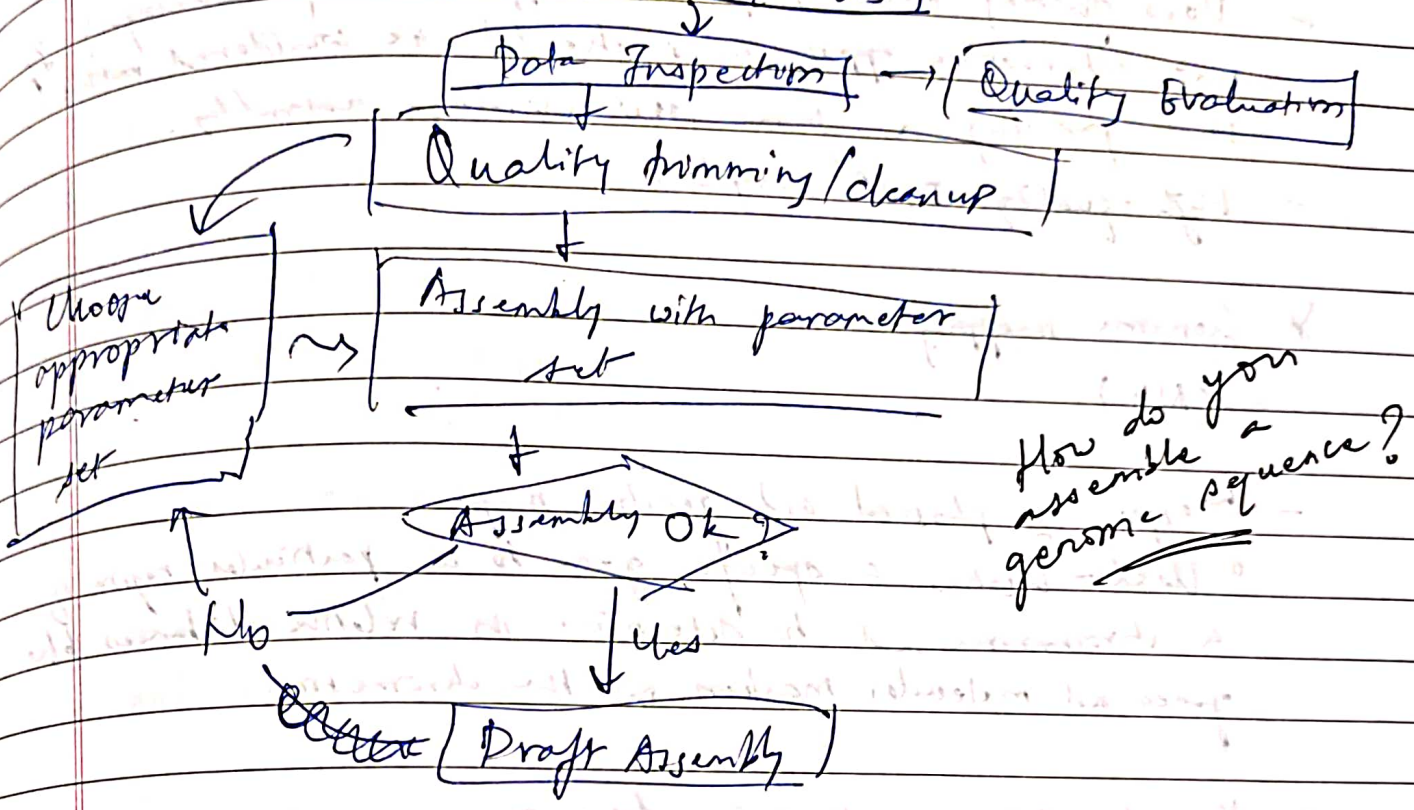
* NGS -
(NKG)

- NGS is a powerful platform that has enabled the sequencing of ^{DNA} ~~1000s~~ ^{billions} to millions of DNA molecules simultaneously.
- This powerful tool is revolutionizing fields such as personalized medicine, genetic diseases and clinical diagnosis by offering a high throughput option with the capability to sequence multiple individuals at the same time.
- have gained an understanding of the underlying cause of disease.

* Assembly of Data from Genome Sequencing -
(NKG)

- In Bioinformatics, ^{assembly} ~~sequencing~~ refers to aligning and merging fragments from a larger DNA seq. in order to reconstruct the original seq. Typically, the short fragments called reads, result from shotgun sequencing of genomic DNA or gene transcript (EST).

- De Novo Gene Assemblies assume no prior knowledge of the source DNA seq. length, layout, or composition.
(Raw read sequences)



Steps - (NK - ①, ②, ③)

- ④ Choose the best sequencing platforms and library ^{preparations.}
- ⑤ Select the best possible DNA source and DNA extraction _{method.}
- ⑥ Check the computational resources and requirements.
- NK - ⑦, ⑧, ⑨ ⑧ Assemble the genome.
- ⑩ Check the assembly quality before annotation.
- ⑪ Genome annotation
- ⑫ Build a searchable and shareable Output format.
- ⑬ Reach out to the community ~~and~~ to refine the assembly and annotation.

* Conclusions -
(NK)
Beginners & small reser

↓ Imp points to consider -
(MK)

- Does funding allow to produce sufficient seq. coverage? If not, alternative approaches should be considered rather than producing a poor low coverage assembly.
- High-quality DNA sample.

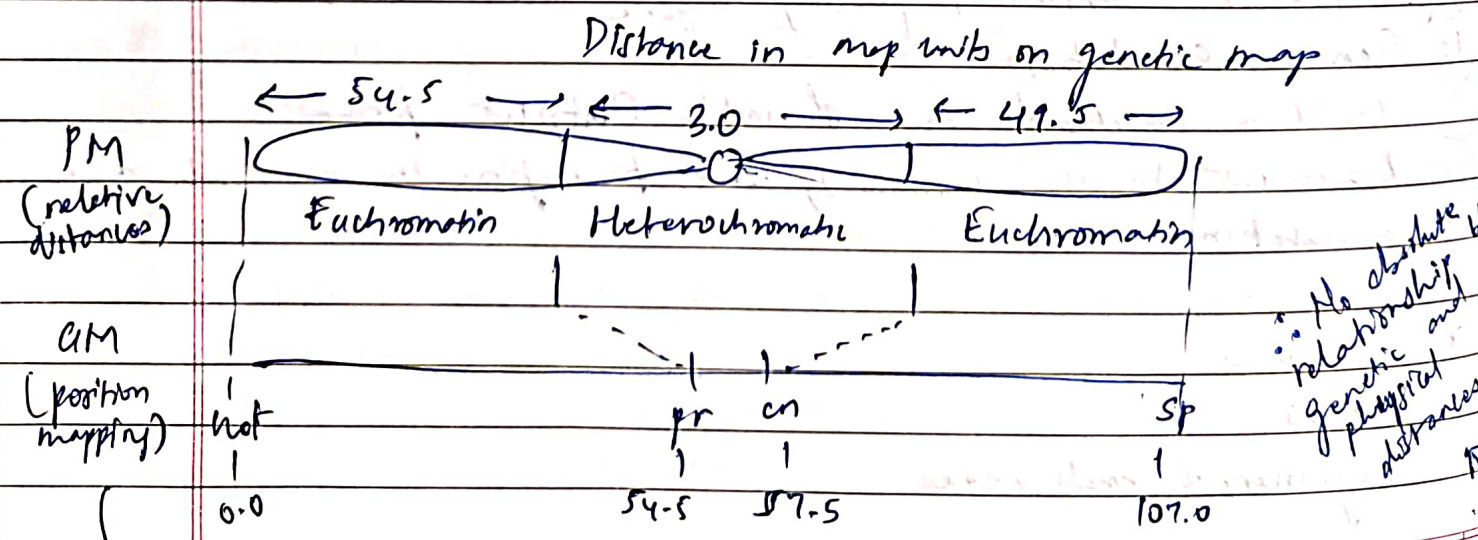
↓ Genome mapping -
(MK)

- Mapping of physical and genetic maps
 - o Used - locating a specific gene to a particular region of a chromosome and to determine its relative distances b/w genes and molecular markers on the chromosome.

↓ How is genomic mapping done -
(MK)

To do this, the genome is first broken up into fragments. The fragments are then replicated up to 10 times in bacterial cells to create a library of DNA clones.

- Genetic vs Physical Distance - Diagram



data may have limited resolution or limited accuracy (recombinational "hotspots" or "coldspots")

Comparing GM and PM -

Saccharomyces cerevisiae chromosome III.

- o physical map obtained by DNA sequencing.
- o the order of the upper 2 markers (gkl1 and dal1) is incorrect on the GM.

Approach for Genom Analysis -

- o CGS, shot gun, DNA markers (RFLP, SSLP, SNP), PCR;
- o microarray (Protein MA, DNA MA);
- o Restriction mapping, FISH, STS, EST.

PMID - 1770 3.239

doi - 10.1038/nrg2144

Identification of Disease Genes 4 Role of Birnbo - eg - Database -

- Green - regular expression genes (no mutations)
- Red - mutated
- Black
- Yellow

Gene Prediction - Continued

FP = false - positive.

* Prediction of Promoter and Regulatory Elements - (PRE)

(MK)

- i.e. these DNA elements directly regulate gene expression.
- PRE are traditionally determined by experimental analysis, but time-consuming and laborious.
- ∴ Computational approach to solve the problem and enhance the potential.

- Problems -

(MK)

Lessons -

- ③ lack of data.
- ④ PRE cannot be translated into protein sequences to ↑ the sensitivity for their detection.
- ⑤ Normally elements are short (6-8 nucleotide) and found in any seq. by random chance ∴ ↑ rate of FP results because of which sensitivity drops and specificity is hampered.

- Solⁿ -

- ④ Predictions based on evolutionary info.
- ⑤ Experimental data.

* PRE in prokaryotes -

- ① Transcription - RNA polymerase
 -70 → recognizes specific sequences upstream of a gene and allows the rest of the enzyme complex to bind.
 -35 and -10 boxes

consensus seq -

(-35) L TATAAT
 TTGACA

(-10)

② Operon

- Transcription Factors (TF) - bind to specific DNA seq. do either enhance or inhibit the function of the RNA polymerase.
- Regulatory elements - specific DNA seqs to which TF bind. (→ may bind in the vicinity of the promoter or bind to a site several 100 bases away from the promoter.)
- Regulatory proteins binding @ long distance can still exert their effect is because of the flexible structure of DNA which is able to bend and exert its effect by bringing the transcription factors in close contact with the RNA polymerase complex.

* PRE in Eukaryotes -

- ① RNA polymerase I - transcription of rRNA and tRNA respectively
- ② II - transcribing protein-encoding genes (or synthesis of mRNAs)
- ③ III - same (same as I)

- Each of them transcribes diff. sets of genes.
- Unlike prokaryotes,

RNA pol II does not directly bind to the promoter; but relies on a dozen or more transcription factors to recognize and bind to the promoter in a specific order before its own...

(NK)

Many genes have a unique initiator sequence (Inr) which is a pyrimidine rich seq with a consensus (C/T)(C/T)CA (C/T)(C/T)

This site coincides with the transcription start site. Mostly TF binding sites are located within 500bp upstream of the Transcription Start site.

Some regulatory sites can be found 10 to 1000 bp away from the gene start site.

RE are located downstream instead of upstream of the transcription start site.

* Prediction Algorithms - (NK)

③ Expression profile based - using profiles constructed from a no. of coexpressed gene seq from the same organism.

④ Phylogenetic footprinting.

- RNA pol II describes the eukaryotic mRNA

① Ab-initio based - (NK)

Conventional Method - Detection of element via matching a consensus seq. pattern represented by regular expression or matching a position - specific scoring matrix (PSSM).

- But in other case, the consensus seqs or the matrices are relatively short, covering 6-10 bases and generates \uparrow rate of FT as a result.
- Solⁿ - New algorithm based on NM, ANM, HMM, MLT. - - -

(a) * Prediction for Prokaryotes -
(CNK)

- But set of simple rules (Wang et al) is accurate which relies on 2 kinds of info -
 (i) gene orientation and intergenic distances of a pair of genes of interest and conserved linkage of the genes based on comparative genomic analysis i.e., gene linkages are focused.
- A scoring scheme is developed to assign spacers with diff. levels of confidence.

CEAMS - Accurate identification of an splicer str., which in turn facilitates the promoter prediction.

- This software scheme not yet available as a computer program thus can be done manually using the rules, but few programs doesn't accept this approach.
- Scoring criteria for splicer prediction
- Prediction Tool for prokaryotes - BPRM.

Prediction for Eukaryotes - (ANK)

- Also because of the ↑ variability of T_A Binding sites, the simple seq. matching often misses true promoter sites, creating F₁.

∴ to ↑ specificity of prediction unique feature of eukaryotic promoter is employed, which is the presence of CpG islands

Many vertebrate genes are characterized by a ↑ density of CpG dinucleotides near the promoter region overlapping the TSS.

∴ CpG island identification can lead to trace promoter in the immediate upstream

Prediction Tool - TSSW

(c) Phylogenetic Footprinting - Based Method - (NH)

- Careats / Red flags are -

① Organisms selected are too closely related (human & chimpanzee) the seq. diff. b/w them may not be sufficient to filter out functional elements.

② Organism evolutionary distances are too long (such as human & psh), long evolutionary divergence may render promoter and other elements undetectable.

③ To extract non-coding seqs upstream of corresponding genes and focus the comparison to this region only, which helps to prevent FP.

④ Predictive value depends on the quality of the subsequent seq. alignment. ∴ Advance alignment program is used.

- ∴ software programs are specifically designed to take advantage of the presence of phylogenetic patterns to make comparisons among a no. of related species to identify putative TFBS.
- Adv - No hairy of probabilistic models reqd.

③ Expression Profiling Based Method -
Xin Xisong

DNA Microarray Introduction

④ Affinity chips - Xin Xisong 26^A

* Basic Concepts

- Prokaryotic v/s Eukaryotic cells
- Cell Organelles, Chromosomes
- DNA - genetic maps
- RNA
- Central Dogma

- The Goal - (NK)

One approach - what happens to the

Large scale approaches - (NK)

↳ Microarrays allow massive parallel measurements in one experiment

- The Southern Blot -

- Microarray - (NK) -

↳ Array - Ordered arrangement of large gene samples.

Background -

↳ Functional genomics - (NK)

o study of obtaining an overall picture of genome functions, including the expression profiles at the mRNA level and the protein level.

(NK)

- o Integrative biology & systems biology studies to understand health & disease states.
- o Drug discovery process.

↳ Gene expression - NK

↳ Microarray -

o Tools used to measure the presence & relative abundance of gene expression in tissues.

- Outline of Microarray -
(NK)

- (I) Interpretation of result
- a. Statistical analysis --

* What is Microarray -
(NK)

- In particular, the amount of mRNA for each gene in a given sample (or a pair of samples) is measured -

- o Parallel
- o High-throughput
- o Large scale
- o Genomic scale

- Historical perspective
(NK)

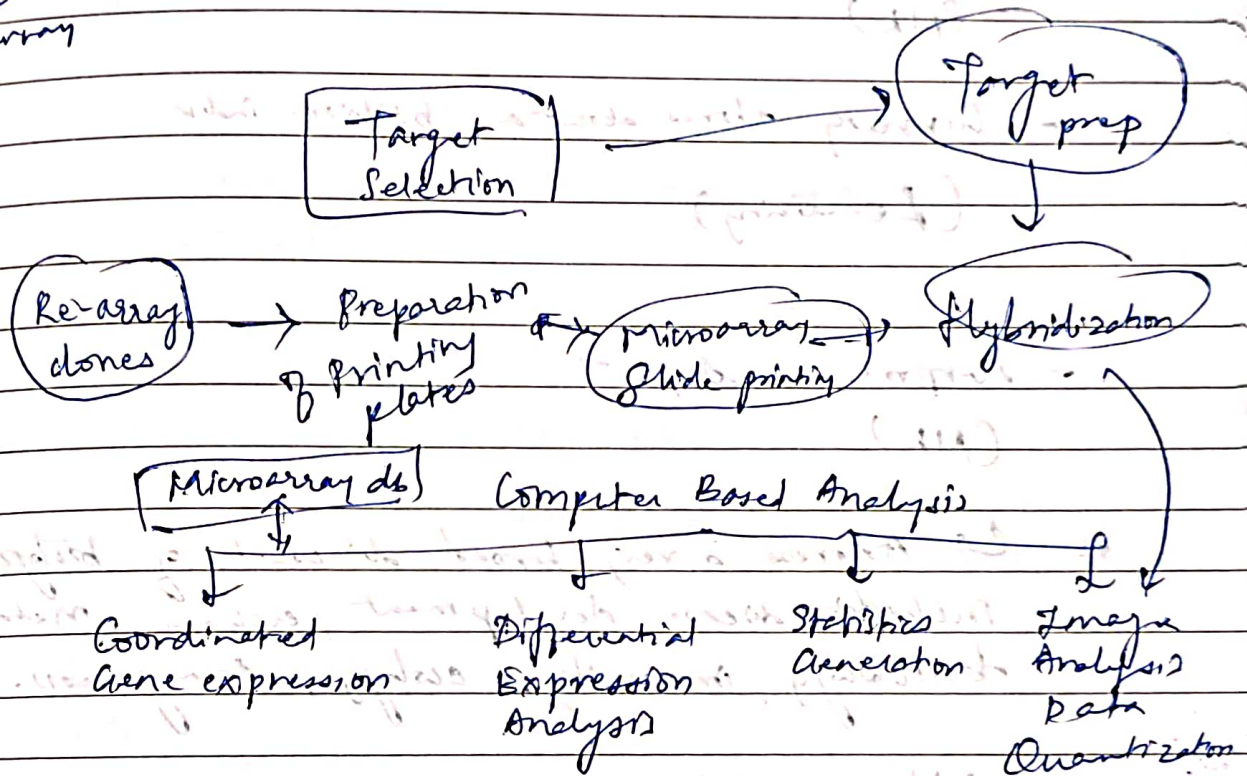
- o Finny sample

GEO [Gene Expression Omnibus]

Also called "Microarray Databases"

add in equal amounts

Hybridize probe to microarray → Scan



* Output of Microarray -

(NK)

- No immediate observations can be made concerning gene expression from raw data.
- Statistical analysis applications are used to interrogate the data for information on gene expression patterns.

* GEO

(MK)

- Convenient for deposition of gene expression data, as required by funding agencies & journals.
- Curated, online resource for gene expression data which is required for analysis.

(MK)

- Currently stores about a billion indiv.

(Remaining)

- Purpose and scope -

(MK)

- * → Address a very broad diversity of biological themes including disease, development, evolution, metabolites, toxicology, immunity, ecology, transgenetics.

↳ Not intended to be used as a Laboratory Information Management System (LIMS) or a pre/post analysis

- 3 main facts of GEO -

(1) MK

(2) Offer simple submission procedures and formats that support complete and well-annotated data deposits from the research community (submission guide).

(3) Provide user-friendly mechanisms that allow users to query, locate, review and download studies

GEO structure / architecture

4 kinds of data records.
(MK)

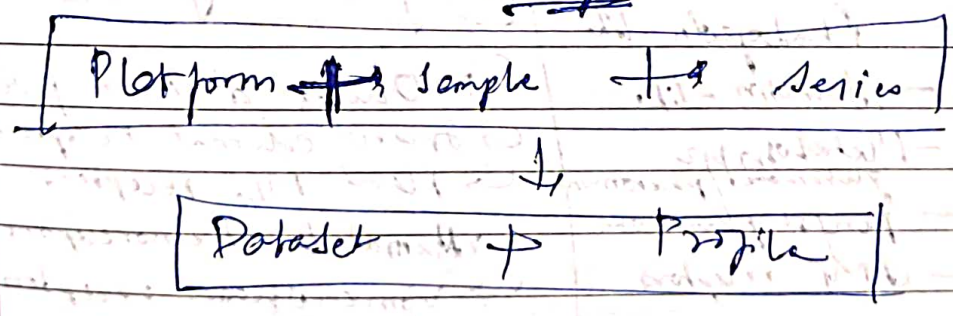
③ Series (GSE) = defines a set of samples and how they are related.

④ Datasets (GDS) = sample data collections assembled by ^{steps} GEO.

↳ Submitters may provide raw data.

Submitted by manufacturer	Submitted by Experimentalists		Curated by NCBI
<u>GPL</u>	<u>GSM</u>	<u>GSE</u>	<u>GDS</u>
Platform Description	Raw / processed spot intensities from a single slide / chip	Grouping of slide / chip data " a single experiment	Grouping of experiments

Data Organization



* Navigating GEO -

G-protein coupled Receptor

(MK)

- Belongs to a family of proteins that act like a molecular switch inside the cells and are involved in transmitting the signal from a variety of stimuli outside a cell to the inside.

(MK)

- When an external signaling molecule binds to a GPCR, it causes a conformational change in the GPCR.
- This change then triggers the interaction b/w the GPCR and a nearby protein.

* Structure -

(MK)

* Types of GPCR -

(2) Frizzled / Smoothened family

(3) Potamine families

- (1) Class A - Rhodopsin-like
 B - Secretin-like
 C - Metabotropic glutamate/pheromone
 D - Fungal pheromone
 E - AMP receptors

- ↳ Olfactory adenylyl cyclase proteins
- ↳ Insect odorant receptors
- ↳ Plant Mlo receptors
- ↳ Nematode chemoreceptors
- ↳ Vomeronasal receptors (V1R and V2R)
- ↳ Taste receptors (T2R)

Main classes based on sequence homology and function.

(4) Orphans
 ↳ Putative/unclassified GPCRs

(5) Non-GPCR families
 ↳ Archaeal / Bacterial / fungal opsins

* Types of Ligands

① Agonists

Ligands which shift the equilibrium in favor of the active state.

② Inverse Agonists

- Ligands which shift the equilibrium in favor of inactive states

③ Neutral Antagonists

- Ligands which do not affect the equilibrium.

* Role of GPCR - (in body)

- ① Hormones
- ② Odorants
- ③ Tastes

* GPCR function -

① Visual sense - Spoons use a photoisomerization rxn to translate electromagnetic radiation into cellular signals. Rhodopsin, for eg, uses the conversion of 11-cis-retinal to all-trans-retinal for this purpose.

② Sense of smell - Receptors of the olfactory epithelium bind odorants (olfactory receptors) and pheromones (vomeronasal receptors).

③ Behavioral & Mood regulation - Receptors in the mammalian brain bind several diff. neurotransmitters, including serotonin and dopamine.

④ Regulation of immune system activity and inflammation -

Chemokine receptors bind ligands that mediate intercellular communication b/w the immune system; receptors such as histamine receptors bind inflammatory mediators & engage target cell types in the inflammatory response.

⑤ Autonomic Nervous System - Both the sympathetic and parasympathetic nervous systems are regulated by GPCRs pathways. These systems are responsible for control of many autonomic functions body such as blood pressure, heart rate, digestive...

* Orphan GPCR [Orphan GPCR] -

- GPCR belongs to a supergene family and \therefore would share seq. similarities.
- The no. of GPCR is about 800, of which more than half are \odot orphans GPCRs.
- The discovery of new GPCRs found by homology screening suffers from one obvious problem, the receptors found lack their pharmacological identities, their natural ligands.
- Without knowing what their endogenous ligands are, we lack the information to understand their physiological role and hence cannot...
- Traditionally the ligand was found first and served to characterize the receptors pharmacologically.
- But by the end of the 1980s, study for identification of GPCRs was obstructed by using homology screening approaches.
- BUT
- GPCRs found through homology approaches lack their natural ligands, thus called as "Orphan" GPCRs.
- Searching for ligands of orphan GPCRs has given birth to the "REVERSE PHARMACOLOGY" approach.
- \therefore Orphan as targets to identify their endogenous ligands started.

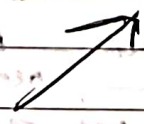
Orphan receptor

Adopted orphan

Similar structure to other identified receptors but whose endogenous ligands has not yet been identified.

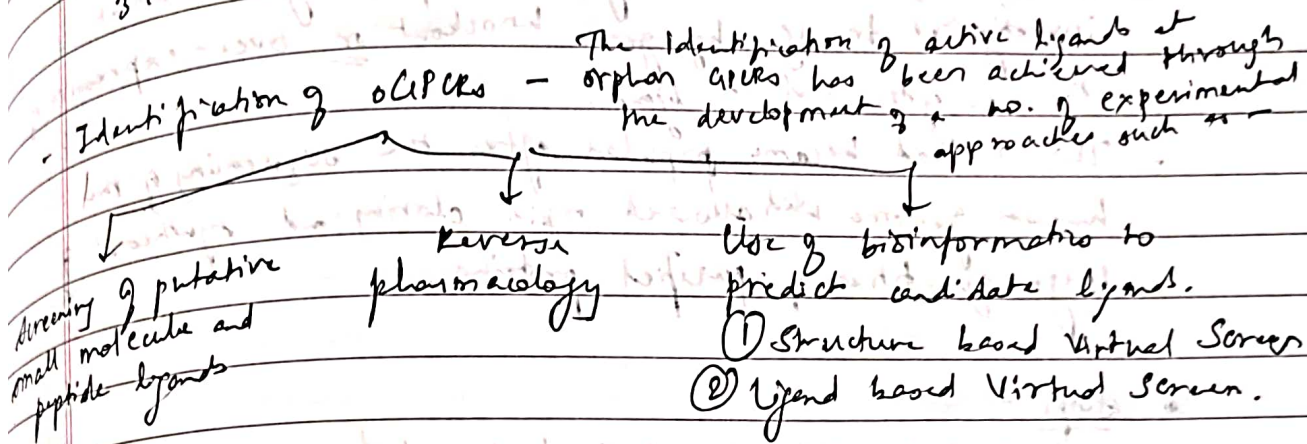
Identifying the endogenous ligand for the orphan receptor.

If a ligand for an orphan receptor is later discovered, the receptor is referred to as an "adopted orphan".



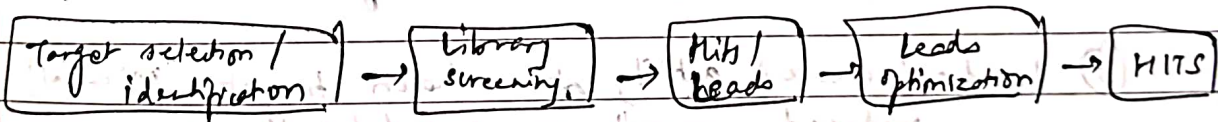
- Classification of known and orphan (non-olfactory) GPCRs -
 - 39% Class I, knowns (Rhodopsin class)
 - 6% Class II

37% unknowns.



* Classical pharmacology / forward pharmacology - (Nk)

- Using the performance of medicinal chemistry, the potency, selectivity and other possessions of these screening hits are optimized to create candidate drugs.
- Compounds are screened in cellular or animal replica of infection to categorize compounds that cause a desirable transformation in phenotype.



(Nk)

- ... disease modifying and screen for compounds that modulate the activity of this purified target.
- Later on compounds are tested in animals to see if they have the desired effect.
- This approach is known as "reverse pharmacology" or "target based drug discovery" (Tdd).

* Reverse Pharmacology - (NK)

- The identification of a receptor sequence is followed by the discovery of the corresponding ligand.
- Subsequently the pharmacological and physiological context is investigated for eg. by gene knockout or over-expression of the receptor and its ligand.
- This method became popular after the sequencing of the human genome which allowed rapid cloning and synthesis of large quantities of purified proteins.

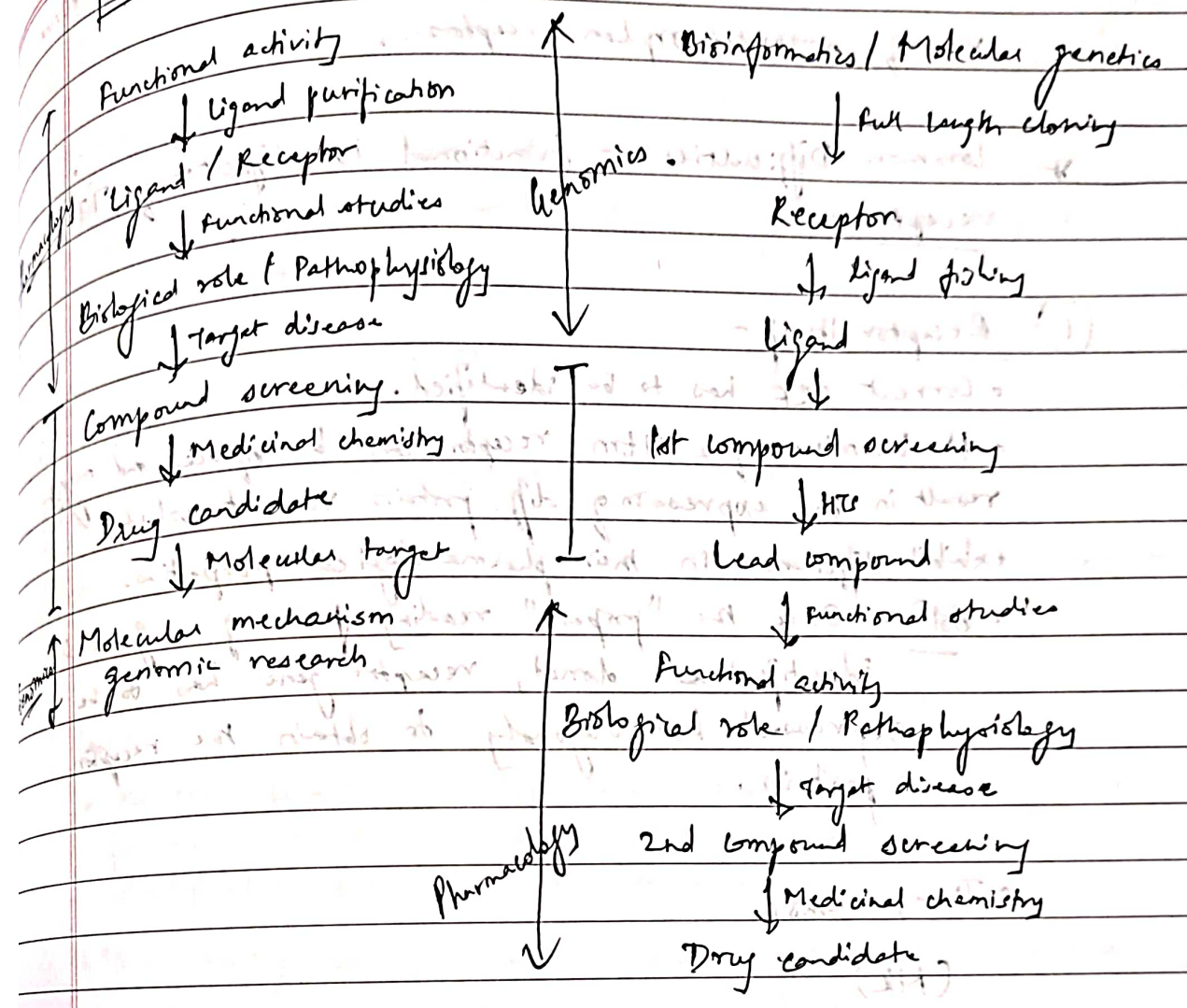
* Scope -

- Understand the MOA at multiple levels of biological organization.
- To optimize safety, efficacy and acceptability of the leads in natural products, based on relevant science.

* Keywords -

- ① Drug Target - Cellular/molecular structures involved in the pathophysiology of interest where... (NK)
- ② Hit - potential activity at a chosen target.
- ③ Lead - increased activity at a chosen target. (potential)
reduced ——— " ——— unrelated target.

Classical approach Reverse pharmacology



⚡ Reverse Pharmacology Methods - (Ⓢ Diagram for all 3 methods) (Colorpharization)

① Library based approach -

② Tissue extract based -

- (a) Isolating active compounds from crude tissue extract at laboratory scale.
- (b) Identifies hits in fractionated tissue samples.

③ Information based -

- (a) Identifies prospective ligands by database screening and testing them on a small no. of selected orphan receptors.

(b) The available information regarding affected cells or tissues is used to ...

no. of suspected orphan receptors ...

Common Difficulties - functional identification of GPCR receptor -

(i) Receptor itself -

- o Correct ORF has to be identified.
- o N-terminus of a GPCR receptor can be spliced, and might result in the expression of diff. protein variants which exhibit differences in their pharmacological properties.
- o Δ ORF - Once the "proper" reading frame of a gene is identified and cloned, receptor gene has to be expressed heterologously to obtain the receptor protein.

o Thus, recomb.

(NHK)

For eg -

(a) The receptor can be tagged with either a short peptide (on the N-terminus).

(b) ^{however} both of the techniques could potentially influence receptor protein targeting, disturb ligand binding, or alter the signaling behavior of the receptor.

(c) Hence, the cell surface expression of the receptor has to be controlled in functional assays.

(NK) (c) the over-expressed receptor might also attract G -Proteins from other endogenously expressed receptors, thus decreasing signals obtained from these receptors.

(d) None of the techniques described is straightforward; all require a considerable amount of work, yet, still they do not guarantee success.

(2) (NK) Related to the receptor protein might arise from its possible interaction with

(3) coupling of the receptor to effectors inside the cell.

(4) The difficulty related to receptor signalling & assay systems are shared by all 3 approaches.

(NK)

However, because of the high demands placed on the assay used for the initial characterization, techniques successfully used for identification purposes are not as numerous.

These assay techniques have to display high reliability and quality, a property that is defined by a high Z-factor.

Summary

(2) The tissue-extract-based approach -

Requires an ~~extra~~ extraordinarily sensitive assay procedure...

Case Study - GPCR

Reverse Pharmacology of orphan GPCRs

(NK)

Later, when high-throughput screening technology was applied to reverse pharmacology, dozens of orphan GPCRs became deorphanized.

Furthermore, novel neuropeptides were discovered.

- Artemisinin -

(NK)

- Artemisinin and its derivatives offer promise as a new class of ^{antimalarial} drugs.
- Best case for reverse pharmacology approach.
- Anti-cancer, asthma and anti-oxidant properties.

(NK)

• Although molecular biological and bioinformatics techniques made the identification of orphan GPCRs responsive, the search for their endogenous ligands has been a challenge.

- This search has given birth to the reverse pharmacology approach, which uses orphan GPCRs as targets to identify endogenous ligands.
- This approach was very successful and has led to over 2 decades to the deorphanization of about 300 GPCRs.

- Searching for Novel Transmitters -

(NK)

- These transmitters are mostly small molecules although few of them are larger polypeptides and they include, neuropeptides, lipid mediators, nucleotides, amino acid and derivatives etc.
- So how could one find which one of these neurotransmitters may bind on orphan GPCRs?

(NK)

• (NK)

- Because its phylogenetic analysis classified ~~GPCR~~ opioid receptor like 1 as a peptidergic (chemical which functions to directly modulate the peptides of the body or brain). GPCR and because the opioid receptor-like is expressed in CNS, peptidergic brain tissue extracts were prepared, purified and fractionated.

(NK)

- Its structural similarities to the opioid peptides made it an immediate star, yet it has been proven not to bind the opioid receptors.
- Further studies need to be carried out.

Three Les of Reverse Pharmacology -

◦ (NK)

- Majority of the drug discoveries would not have been developed or their development would have been delayed significantly in the absence of the scientific or technical contributions from the pharmaceutical companies.

◦ (NK)

- Best of public and private sector partners comprising academia and industry should come together to reap significant benefits from these seemingly low fashionable but highly gifted explorations based on traditional knowledge.

◦ (NK)

- Reverse Pharmacology approaches need to be developed further and optimized as novel means for fast track drug discovery and development of newer, safer and effective drugs.

Introduction to Proteomics

PAGE No.

DATE

* Protein Identification with Ab -

(NK)

- Each tip \rightarrow of the "Y" of an Ab contains 2 paratope (analogous to a lock). This is specific for one particular epitope (analogous to a key) on an antigen.

- This architecture allowing these 2 structures to bind together with precision.

- Using this binding mechanism, an Ab can tag a microbe or an infected cell for attack by other parts of the immune system, or can neutralize it directly. (for eg, by blocking a part of a virus that is essential for its invasion).

- Abs can be used to visualize the location of specific proteins within the cell.

* The benefit of detecting Ab to flow -

(MK)

- These methods used is -

- ① Immunocytochemistry / Immunofluorescence - refers to the visualization of specific antigens in culture cells.
- ② Immunohistochemistry - refers to their visualization in prepared tissue sections.

- Immunofluorescence and immunohistochemistry - use antibodies for detection and localization of proteins and other antigens within biological ~~systems~~ samples.

(NK)

Common steps -

① fixation of the sample onto slides or plates, blocking in order to avoid unspecific signal, detection by the use of one or two antibodies, and analysis.

② (NK).

③ this is an imp. step because immune cells may be ~~present~~ present within the tissue.

④ such approach avoids an unspecific response because primary Abs can bind to antibody receptors present in immune cells membrane. (eg - Mast cell).

* Immunofluorescence (IF) -

④ Diagram

(NK)

- 1st - one Ab is used which binds directly to the target Ag, ∴ this technique is called DIF

- 2nd: ~~the~~ it is called IIF when 2 Ab are used & because the Ab allowing for detection binds to a 1st Ab that recognizes the ~~Ab~~ antigen on that slide.

* Conclusion -

- ∴ Necessary to -

① permeabilizing reagents which expose the epitope and highlight the localization of antigen because there are some antibodies which cannot go into the nucleus; ∴ the epitope will not be achieved and no signal will be produced.

- ② A sample incubated only with the 2^o Ab to determine non-specific binding sites.
- ③ Control slide containing cells that either do not express ad/or...

* Immunohistochemistry -

(IHC)

- ... precipitate at the location of the protein or fluorescent detection, in which a fluorophore is ^{conjugated} ~~combined~~ with the Ab and can be visualized using fluorescence microscopy.
- Can be performed using samples with different previous treatments such as ^{cryo} preserved slides or paraffin-embedded tissue slides.

(IHC)

- ④ Negative: controls could be slides previously exposed to Ab-specific serum (e.g. rabbit's serum if rabbit's Ab is used) or isotype-specific Ig as primary antibody.

⑤ Controls ^{are} also useful for determination of non-specific binding sites to secondary Ab and optimal primary Ab dilution.

⑥ Positive controls should also be included to validate proper activity of the d.f. reagents.

⑦ ^{IHC} can be done also with fluorochrome labeled antibodies.

* Protein Identification -

- ① Edman's Degradation (ED) -
(IHC)

- To solve the problem of damaging the protein by hydrolyzing conditions, Peter Edman created a new way of labeling and cleaving the peptide. Edman thought of a way of removing only one residue at a time, which did not damage the overall sequencing.

Procedure -

- ① Add Phenyl isothiocyanate, which creates a phenylthio carbamoyl derivative with the N-terminal.
- ② The N-terminal is then cleaved under less harsh acidic conditions, creating a cyclic compound of phenylthio hydantoin PTH-amino acid.
- ③ This does not damage the protein and leaves 2 constituents of the peptide.
- ④ This method can be repeated for the rest of the residues, separating one residue at a time.

(CNK) This can be done by denaturing the protein and heating it and adding HCl for a long time.

- This causes the individual amino acids to be separated, and they can be separated by ion exchange chromatography.
- They are then dyed with ninhydrin and the amount of amino acid can be determined by the amount of optical absorbance.
- This way, the composition but not the sequence can be determined.

* Proteomics -

(CNK)

- The separated peptides can be isolated by chromatography.
- They can be sequenced using the Edman method, because of their smaller size.
- In order to put together all the sequences of the diff. peptides, a method of overlapping peptides is used.
- The strategy of divide and conquer ...

10.30
11-12

$$x^2 + y^2 = \text{circle}$$

(NK)

- However, this method is limited in analyzing larger sized proteins (more than 100 a.a.) because of secondary hydrogen bond interference.
- Other weak intermolecular bonding such as hydrophobic interactions cannot be properly predicted.
- Only the linear seq. of a protein can be properly predicted assuming the seq. is small enough.

* Shotgun proteomics for Proteomic Profile - (Self)

- Assignment - EXPASY

* Protein Microarray -