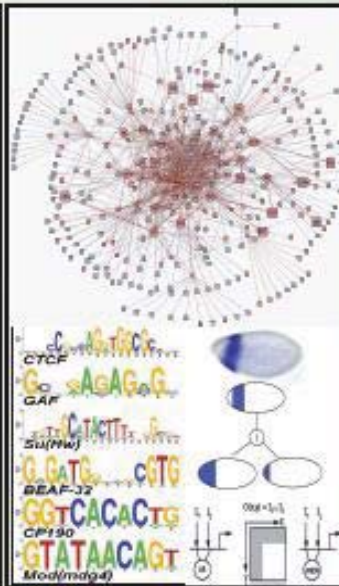# 6.047 (U) - 6.878 (G, TQE) - HST.507

Fall 2015

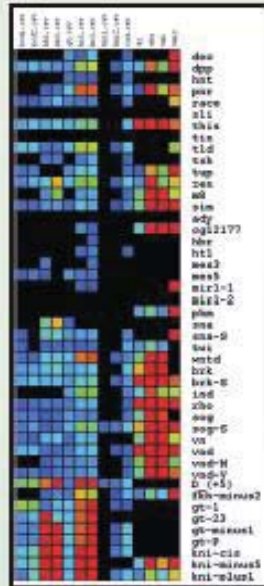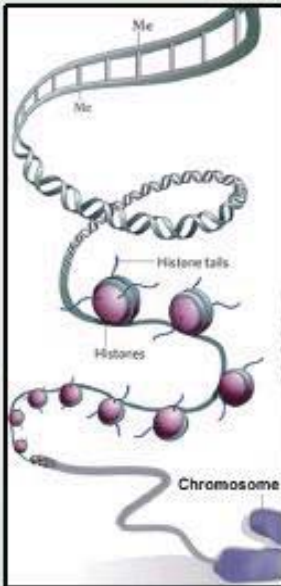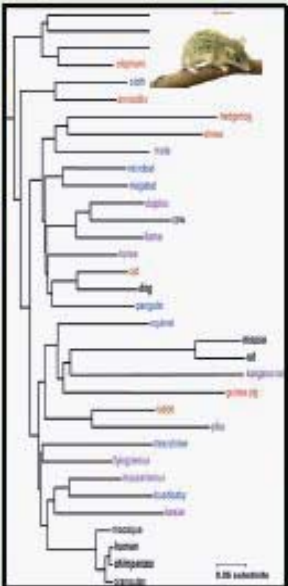# Computational Biology
## Genomes - Networks - Evolution

Learn about:



Comparative genomics — Epigenomics — Functional genomics — Motifs & networks — Phylogenomics — Personal genomics

… and much much more

GWAS

## Prof. Manolis Kellis - MIT / CSAIL / Broad Institute

Covers the algorithmic and machine learning foundations of computational biology combining theory with practice. We cover both foundational topics in computational biology, and current research frontiers. We study fundamental techniques, recent advances in the field, and work directly with current large-scale biological datasets.

# Computational Biology: Genomes, Networks, Evolution

**Rapid database search**

## MIT 6.047 / 6.878
## HSPH IMI.231
## HST.507

**Protein interaction network**

## Prof. Manolis Kellis



**Genome duplication**

# I. Administrivia

Introduction to the course and its goals

Course organization and content

Homework and Quiz

Term Project

# Introductions

- **Lecturer**
  - Manolis Kellis
    (MIT CSAIL, Computational Biology, Broad Institute)
  - My own research:
    Comparative genomics, Gene Regulation, Evolution,
    Epigenomics, Phylogenomics, etc

# Course Information

- Lectures
  - TR 1pm – 2:30
- Recitations:
  - On Friday at 3pm
  - Recitations at MIT (HST/HSPH students can join)
  - All handouts, lectures, notes, etc will be posted here.

- Course calendar:
  - On Google, add public calendar: "6.047 Lectures" and "6.047 due dates"

# Goals for the term

- ## Introduction to computational biology
  - Fundamental problems in computational biology
  - Algorithmic/machine learning techniques for data analysis
  - Research directions for active participation in the field
  - Understanding *how* methods work

- ## Ability to tackle research
  - Problem set questions: algorithmic rigorous thinking
  - Programming assignments:
    → hands-on experience w/ real datasets
  - Final project experience:
    → propose and carry out independent original research
    → present findings in conference format (written, oral)

# Course content

# Computation & Biology | Foundations & Frontiers

- Duality #1 (x-axis):  Computation and Biology
  - **Important, relevant, current biology**:
    - → Important biological problems
  - **Fundamental computer science**:
    - → General techniques, principles
- Duality #2 (y-axis): Foundations and Frontiers
  - **Foundations**:
  - well-defined problems, general methodologies
  - 'The classics' of the field
  - **Frontiers:**
  - in-depth look at complex, current problems, open questions
  - combine techniques learned
  - opens to projects, research directions

# Course organized around bio/comp modules

- Each module corresponds to an active area of research
  - **1: Comparative genomics: Align/model genomes, DP, HMMs**
  - **2: Genes and Transcripts: RNA-seq, clustering, structure**
  - **3: Regulation: Epigenomics, TFs, Motifs, Network inference**
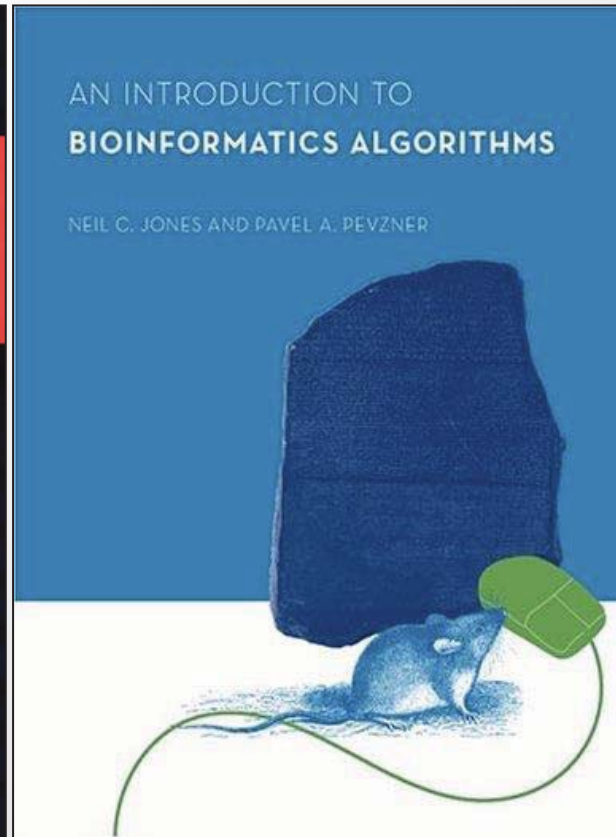  - **4: Variation: Genetics, Human history, heritability, eQTLs**
  - **5: Evolution: Phylogeny, evolutionary sigs, WGD, assembly**
  - **6: Frontiers: Personal/Disease, 3D genomes, Pharma, Synth**
- For each module:  First half ⇔ the foundations
  - Dynamic programming, string matching, hashing, HMMs, EM, Gibbs Sampling, Clustering, Classification, Feature selection, SVMs, CRFs, Context-Free Grammars, phylogenetics, gene / species trees, evolutionary models, GWAS, disease mapping
- For each module:  Second half ⇔ the frontiers
  - Evolutionary signatures, Transcript analysis, lincRNAs, Network inference and analysis, Epigenomics, Recent human selection and ancestry, chromatin regulation, Missing heritability, 3D
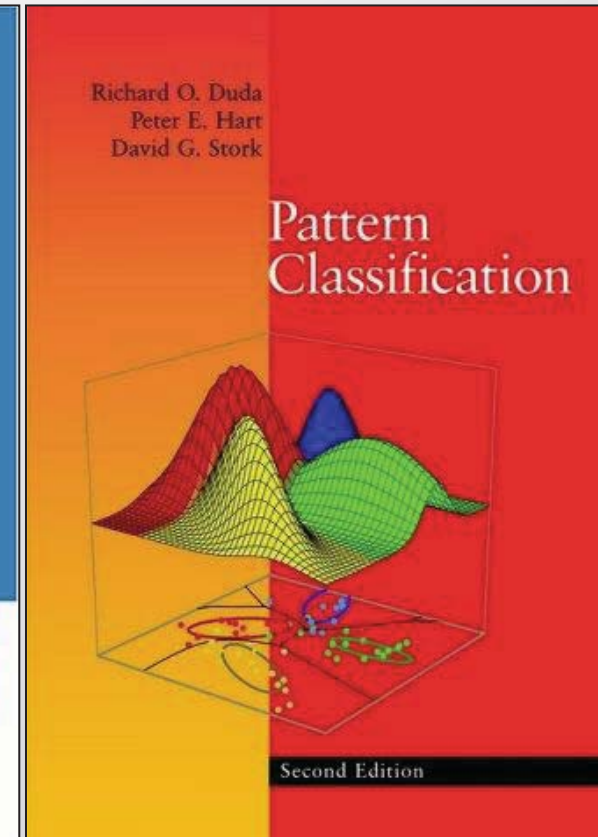
| Project | Psets | Week | Date | Topic | | Lec | Topic | Read* |
|---|---|---|---|---|---|---|---|---|
| Describe your previous research, areas of interest in computational biology, type of project that best fits your interests. Post in a profile that lets your classmates know you and find potential partners. **Project profile due Tue 9/29** | PS1 out on:L1-L5 | 1 | Thu, Sep 10 | **Introduction** | | L1 | Intro: Biology, Algorithms, Machine Learning, Course Overview | 1 |
| | | | Fri, Sep 11 | | | R1 | Recitation 1: Biology and Probability Review | |
| | | 2 | Tue, Sep 15 | Module I: Aligning and Modeling Genomes | Foundations | L2 | Alignment I: Dynamic Programming, Global and local alignment | 2 |
| | | | Thu, Sep 17 | | | L3 | Alignment II: Database search, Rapid string matching, BLAST, BLOSUM | 3 |
| | | | Fri, Sep 18 | | | R2 | Recitation 2: Deriving Parameters of Alignment, Multiple Alignment | |
| | | 3 | Tue, Sep 22 | | Frontiers | L4 | Hidden Markov Models Part 1: Evaluation/Parsing, Viterbi, Forward algorithms | 7 |
| | due Tue 9/29 | | Thu, Sep 24 | | | L5 | Hidden Markov Models Part 2: Posterior Decoding, Learning, Baum-Welch | 8 |
| | | | Fri, Sep 25 | | | | No classes - student holiday | |
| | | | Fri, Sep 25 | | | | Project Intro: about the projects, self introductions, mentor intro, example projects, teamwork 32D-507 | |
| Identify previous project proposals, recent papers, and potential partners that match your areas of interest. List initial project ideas and partners. **Project area/team due Tue 10/6** | PS2 out on:L6-R4 | 4 | Tue, Sep 29 | Module II: Gene Expression and Networks | Foundations | L6 | Transcript structure: GenScan, RNA-seq, Mapping, De novo Assembly, Diff Expr | 12.3 |
| | | | Thu, Oct 1 | | | L7 | Expression Analysis: Clustering/Classification, K-means, Hierarchical, Bayesian | 15,16 |
| | | | Fri, Oct 2 | | | R3 | Recitation 3: Affinity Propagation Clustering and Random Forest Classification | |
| Form teams of two, specify project goals, division of work, milestones, datasets, challenges Prepare slide presentation for the class and the mentors. **Project proposal due Tue 10/20. Presented on Fri 10/23** | | 5 | Tue, Oct 6 | | Frontiers | L8 | Networks I: Bayesian inference, deep learning, network dynamics | 20,21 |
| | due Tue 10/13 | | Thu, Oct 8 | | | L9 | Networks II: Network learning, structure, spectral methods | 20,21 |
| | | | Fri, Oct 9 | | | R4 | Recitation 4: Small and Large Regulatory RNAs: lincRNA, miRNA, piRNA… | |
| | | | Fri, Oct 9 | | | | Project Planning: research areas, initial ideas, type of project, mentor matching, finding partners 32D-507 | |
| | PS3 out on:L10-R6 | 6 | Tue, Oct 13 | Module III: Gene Regulation & Epigenomics | Foundations | | No Classes - Monday Schedule - October 13, 2015 | |
| | | | Thu, Oct 15 | | | L10 | Regulatory Motifs: Discovery, Representation, PBMs, Gibbs Sampling, EM | 17 |
| | | | Fri, Oct 16 | | | R5 | Recitation 5: Gapped Motif Discovery, DNAShape, PBMs, Selex | |
| | | 7 | Tue, Oct 20 | | Frontiers | L11 | Epigenomics: ChIP-Seq, Read mapping, Peak calling, IDR, Chromatin states | 19 |
| | due Tue 10/27 | | Thu, Oct 22 | | | L12 | RNA modifications: RNA editing, Translation regulation, Splicing regulation | 11 |
| | | | Fri, Oct 23 | | | R6 | Recitation 6: Dimensionality Reduction | |
| | | | Fri, Oct 23 | | | | Project feedback: Prepare 2-3 slide presentation of your term project for your mentor. 32D-507 at 4-5pm | |
| Evaluate/discuss three peer proposals, NIH review format. **Review Panels Fri 10/30 Reviews back Tue 11/3** | PS4 out on:L13-R8 | 8 | Tue, Oct 27 | Module IV: Population and Disease Genetics | Foundations | L13 | Resolving human ancestry and human history from genetic data | 29 |
| | | | Thu, Oct 29 | | | L14 | Disease Association Mapping, GWAS, organismal phenotypes | 31 |
| | | | Fri, Oct 30 | | | R7 | Recitation 7: Robinson-Foulds Distance and Coalescent Process | |
| | | | Fri, Oct 30 | | | | Panel Discussion: reconciling critiques, strategies for improvement, feedback to author 32D-507 | |
| Address peer evaluations, revise aims, scope, list of final deliverables / goals. **Response due Thu 11/12** | | 9 | Tue, Nov 3 | | Frontiers | L15 | Quantitative trait mapping, molecular traits, eQTLs | 32 |
| | due Tue 11/10 | | Thu, Nov 5 | | | L16 | Missing Heritability, Complex Traits, Interpret GWAS, Rank-based enrichment | 33 |
| | | | Fri, Nov 6 | | | R8 | Recitation 8: Suffix Trees and Arrays | |
| | PS5 out on:L17-R10 | 10 | Tue, Nov 10 | Module V: Comparative Genomics and Evolution | Foundations | | No lecture, veterans day holiday - Monday/Tuesday | |
| | | | Thu, Nov 12 | | | L17 | Comparative genomics and Evolutionary signatures | 4 |
| Continue making substantial progress on proposed milstones. Write outline of final report. **Midcourse report due Thu 11/19. Score projection 11/24** | | | Fri, Nov 13 | | | R9 | Recitation 9: Review of Phylogeny and Molecular Evolution | |
| | due Tue 12/1 | 11 | Tue, Nov 17 | | Frontiers | L18 | Phylogenetics: Molecular evolution, Tree building, Phylogenetic inference | 27 |
| | | | Thu, Nov 19 | | | L19 | Phylogenomics: Gene/species trees, reconciliation, recombination graphs | 28 |
| | | | Fri, Nov 20 | | | R10 | Recitation 10: Linkage Disequilibrium, Haplotype Phasing, and Genotype Imputation | |
| | No more psets! (work on your final project) | 12 | Tue, Nov 24 | Module VI: Current Research Directions | Frontiers | | In Class Quiz (the only quiz - the class has no final exam) - covers L1-R11 | |
| | | | Thu, Nov 26 | | | | No lecture, thanksgiving break - Thu Nov 26, 2015 | |
| Complete your milestones, finalize results, figures, write-up in conference publication format. As par of report, comment on your overall project experience. **Written report due Sun 12/6** | | | Fri, Nov 27 | | | | No recitation, thanksgiving break | |
| | | 13 | Tue, Dec 1 | | | L20 | Personal Genomics, Disease Epigenomics: Systems approaches to disease | 34,36 |
| | | | Thu, Dec 3 | | | L21 | Three-dimensional chromatin interactions: 3C, 5C, HiC, ChIA-Pet | 30 |
| | | | Fri, Dec 4 | | | R11 | Recitation 11: Project Tips - Write-up, Slides, Final Presentation in 32D-507 | |
| Conference format slide pres. **Talks on Thu 12/10** | | 14 | Tue, Dec 8 | | | L22 | Genome Engineering with CRISPR/Cas9 and related technologies | |
| | | | Thu, Dec 10 | | | | Final Presentations - Part I (1pm). 32-141 | |
| | | | Thu, Dec 10 | | | | Final Presentations - Part II (3pm). 32D-507 | |

\* readings refer to chapters in compiled 2014 scribe notes, available in the materials folder

\*\* recitation topics will be adjusted to respond to lecture and student needs

# Textbook / class notes / resources

# (Optional) Books for the Course



**Durbin, Eddy, Krogh, Mitchison**   **Jones, Pevzner**   **Duda, Hart, Stork**

Availability:  BU Coop, amazon.com (~$40-60)

All three books on reserve at the MIT and BU Engineering libraries

# New this year!! Book for the Course

**Computational Biology:
Genomes, Networks, Evolution**

MIT Course 6.047/6.878

**Manolis Kellis & all of you!**

… being compiled this year by students like you!

… actually, including you!

Availability:  Online PDF

# Lectures and Scribing

- Each lecture will have a dedicated scribe who will take notes on the lecture
  - Please sign up to scribe for lecture on the sheet being passed around
- Build on notes from previous years
  - Available on course website
- Complete draft of scribe notes: before prev. lecture
  - Unless it's not there from previous year (this is rare)
- Final draft of scribe notes due 6 days after lecture
  - Your grade depends on the improvement from previous year and completeness
- Some lectures need more work: multiple scribes
- Some tasks are better-suited to you than just scribing
  - E.g. figures, references, layout, macros, let us know!

# Scribing details – DropBox 6047_book LaTex

# Sign up here if you haven't already

| Date | Lecture | Ch | Scribe(s) |
|------|---------|----|-----------|
| Sep. 10 | Intro: Biology, Algorithms, Machine Learning, Course Overview | 1 | Jonathan Li |
| Sep. 15 | Alignment 1: Dynamic Programming, Global and local alignment | 2 | Jesse Tordoff, Thrasyvoulos Karydis |
| Sep. 17 | Alignment 2: Database search, Rapid string matching, BLAST, BLOSUM | 3 | Heather Sweeney, Eric Bartell |
| Sep. 22 | Hidden Markov Models Part 1: Evaluation / Parsing, Viterbi/Forward algorithms | 7 | Anastasiya Belyaeva |
| Sep. 24 | Hidden Markov Models Part 2: Posterior Decoding / Learning Baum Welch | 8 | PH Zhou |
| Sep. 29 | Transcript structure: GenScan, RNA-seq, Mapping, De novo Assembly, Diff Expr | 12.3 | Alex Genshaft |
| Oct. 1 | Expression Analysis: Clustering/Classification, K-means, Hierarchical, Bayesian | 15,16 | Ge Liu |
| Oct. 6 | Networks I: inference, structure spectral analysis | 20,21 | Karthik Murugadoss |
| Oct. 8 | Networks II: Bayesian methods, dynamics, deep learning | 20,21 | Max Shen |
| Oct. 15 | Regulatory Motifs: Discovery, Representation, PBMs, Gibbs Sampling, EM | 17 | ethan evans |
| Oct. 20 | Epigenomics: ChIP-Seq, Read mapping, Peak Calling, IDR, Chromatin States | 19 | Alvin Shi, Connor Duffy |
| Oct. 22 | RNA modifications: RNA editing, translation regulation, splicing regulation | 11 | Narek Dshkhunyan |
| Oct. 27 | Resolving human ancestry and human history from genetic data | 29 | Fernando Varela |
| Oct. 29 | Disease Association Mapping, GWAS, organismal phenotypes | 31 | Sophia Liu, Aurora Alvarez-Buylla |
| Nov. 3 | Quantitative trait mapping, molecular traits, eQTLs | 32 | Giri Anand |
| Nov. 5 | Missing Heritability, Complex Traits, Interpret GWAS, Rank-based enrichment | 33 | Nolan Kamitaki |
| Nov. 12 | Comparative genomics and Evolutionary signatures | 4 | Misha Jamy |
| Nov. 17 | Phylogenetics: Molecular evolution, Tree building, Phylogenetic inference | 27 | Ava Soleimany |
| Nov. 19 | Phylogenomics: Gene/species trees, reconciliation, recombination graphs | 28 | Anne Kim |
| Dec. 1 | Personal Genomics, Disease Epigenomics: Systems approaches to disease | 34,36 | Deniz Aksel, Molly Schmidt |
| Dec. 3 | Three-dimentional chromatin interactions: 3C, 5C, HiC, ChIA-Pet | 30 | Joseph Cunningham |
| Dec. 8 | Genome Engineering with CRISPR/Cas9 and related technologies | | Eunice Wu |

| | Slides | Audio | Notes | Video1 | Video2 | | |
|---|---|---|---|---|---|---|---|
| **Module I: Comparative Genomics** | | | | | | Foundations | **Lecture 1 - Intro and Overview** — Administrivia, Genomes, Inf... |
| | | | | | | | **Lecture 2 - Dynamic Programming / Sequence Alignment** — Dynamic Programming, Sequence Alignment |
| | | | | | | | **Lecture 3 - Hashing, Database Search, BLAST algorithm** — Sequence alignment II, review, local vs. global alignment, semi-numerical string matching, BLAST algorithm, probabilistic interpretation of score matrices (addendum - Linear-time deterministic string matching) |
| | | | | | | Frontiers | **Lecture 4 - Comparative Genomics I - Evolutionary Signatures1** — Evolutionary signatures of protein-coding genes |
| | | | | | | | **Lecture 5 - Comparative Genomics II - Evolutionary Signatures2** — Evolutionary signatures for diverse classes of functional elements |
| | | | | | | | **Lecture 5 - Comparative Genomics III - Evolution** — Mechanisms of evolutionary change, Genome Duplication |
| **Module II: Coding and Non-coding Genes** | | | | | | Foundations | **Lecture 6 - Hidden Markov Models I - Generation, Evaluation, Parsing** — Intro to HMMs |
| | | | | | | | **Lecture 7 - Hidden Markov Models II: Posterior Decoding, Learning** — Increasing state space, Posterior decoding, Supervised/Unsupervised Learning |
| | | | | | | Frontiers | **Lecture 8 - Gene Identification: Gene structure, Semi-Markov, CRFs** — Capturing gene structure, Semi-Markov models, Conditional Random Fields, Emerging lines of evidence |
| | | | | | | | **Lecture 9 - RNA structure** — RNA world, folding algorithms, DP nussinov, energy models, probabilistic models, genomics of ncRNAs |
| **Module III: Networks and Gene Regulation** | | | | | | Foundations | **Lecture 10A - Expression Clustering** — Module III intro, Gene regulation, Microarrays, Expression Clustering, K-means, Fuzzy K-means, Expectation Maximization, Hierarchical Clustering, Hypergeometric |
| | | | | | | | **Lecture 10B - Classification** — Clustering reprise, Bayesian Classification, Naive Bayes, Support Vector Machines |
| | | | | | | | **Lecture 11 - Regulatory Motif Discovery** — TF binding, EM, EM extensions, Gibbs Sampling, Information Content, DNA/protein motifs |
| | | | | | | Frontiers | **Lecture 12 - Regulatory Genomics** — De novo motif discovery using comparative genomics, target prediction and motif instance identification, microRNA hairpin prediction, mature microRNA prediction |
| | | | | | | | **Lecture 13 - Regulatory Networks** — Network structure, network inference, network-based prediction |
| | | | | | | | **Lecture 14 - Epigenomics and chromatin states** — Using combinations of chromatin marks to interpret the human genome |
| **Module IV: Evolution** | | | | | | Foundations | **Lecture 15 - Phylogenetics, Evolutionary Models, Tree Building** — Introduction to phylogenetics, models of evolution, and tree building algorithms |
| | | | | | | | **Lecture 16 - Phylogenomics** — Studying phylogenetics at the genome level, gene/species tree reconciliation, coalescence |
| | | | | | | | **Lecture 17 - Population genomics** — Statistical genetics and human disease mapping |
| | | | | | | Frontiers | **Lecture 18 - Population genetics and recent selection** |
| | | | | | | | **Lecture 19 - Population history** — Population genomics and recent human history |
| **Frontiers** | | | | | | Guest Lectures | **Lecture 20 - Metabolic modeling** — Systems biology for modeling metabolism and regulation |
| | | | | | | | **Lecture 21 - Bacterial Genomics and Microbiomics** — Systems biology for modeling metabolism and regulation |
| | | | | | | | **Lecture 22 - Large intergenic non-coding RNAs** — Genome regulation by large intergenic non-coding RNAs |

17

**Lecture feedback:**

1. Your interest in the overall topic: 1-5
2. The material actually presented 1-5
3. Quality of presentation
   - Quality of slides 1-5
   - Clarify of explanations 1-5
   - Usefulness of lecture notes 1-5
   - Were questions adequately answered 1-5
4. Pace:
   - Difficulty of the material: too easy - just right - too hard
   - Amount of material covered: too little - just right - too much
   - Pace of the lecture: too slow - just right - too fast
5. Comprehension (for each topic)
   - <20%, 20-40%, 40-60%, 60-80%, >80%

# Homeworks and quiz

# Details on Problem sets

- Each problem emphasizes one lecture (or two)
  - Practical problem: gain experience in techniques, write code, download datasets, carry out analysis, interpret your results, learn about behavior of problem/method
  - Theoretical problem: pen/paper, explore algorithmic / statistical / machine learning aspect in detail/depth. (Typically additional advanced problem for 6.878)
- Due Tuesdays at 8pm
  - Late policy: we are flexible, give or take a few hours
  - If more than a few hours, need prior arrangements, extensions typically not granted, except special circ.
- Submit all homeworks online
  - No solutions distributed. If you've solved them, you know what you needed to learn/discover/achieve.

# Details on the in-class quiz

- It's not a midterm, and it's not a final exam
  - It's a quiz, friendly, fun, interesting, cute, fuzzy
- Demonstrate mastery of the material in 4 modules
  - Understand key points emphasized in lecture
  - Understand subtleties revealed in the psets
  - Ability to apply new skills to solve practical problems
- Types of questions
  - Knowledge questions: T/F justify, multiple choice
  - Deeper understanding questions: short answers
  - Practical problems: work through simple algorithm
  - Design problem(s): new/modified algorithm, need both knowledge and new idea, argue correctness

# Final Project

# Final Project: Original Research in Comp Bio

- A major aspect of the course is preparing you for original research in computational biology.
  - Framing a biological problem computationally
  - Gathering relevant literature and datasets
  - Solving it using new algorithms, machine learning
  - Interpreting the results biologically
- Also ability to present your ideas and research
  - Crafting a research proposal (fellowships/grants)
  - Working in teams of complementary skill sets
  - Review peer proposals, find flaws, suggest imprvmts
  - Receiving feedback and revising your proposal
  - Writing up your results in a scientific paper format
  - Presenting a research talk to a scientific audience
- Term project experience mirrors this process

# It's a team project

- Please make an effort to meet your peers!
- Form teams early with complementary expertise

**Final Project at a Glance**

*Project planning* | *Project execution*

| Project | Psets | Week | Date | Topic | |
|---|---|---|---|---|---|
| Describe your previous research, areas of interest in computational biology, type of project that best fits your interests. Post in a profile that lets your classmates know you and find potential partners. **Project profile Due Mon 9/23** with PS1 | PS1 out on:L1-L5<br><br>**due Mon 9/23** | 1 | Thu, Sep 04 | Introduction | |
| | | | Fri, Sep 05 | | |
| | | 2 | Tue, Sep 10 | Module I: Aligning and Modeling Genomes | Foundations |
| | | | Thu, Sep 12 | | |
| | | | Fri, Sep 13 | | |
| | | 3 | Tue, Sep 17 | | Frontiers |
| | | | Thu, Sep 19 | | |
| | | | Fri, Sep 20 | | |
| | | | Fri, Sep 20 | | Project Intro: ★ |
| Identify previous project proposals, recent papers, and potential partners that match your areas of interest. List initial project ideas and partners. **Project area/team Due Mon 10/7** with PS2 | PS2 out on:L6-R5<br><br>**due Mon 10/7** | 4 | Tue, Sep 24 | Module II: Genes and Transcripts | Foundations |
| | | | Thu, Sep 26 | | |
| | | | Fri, Sep 27 | | |
| | | 5 | Tue, Oct 01 | | Frontiers |
| | | | Thu, Oct 03 | | |
| | | | Fri, Oct 04 | | |
| | | | Fri, Oct 04 | | Project Planni ★ |
| Form teams of two, specify project goals, division of work, milestones, datasets, challenges Prepare slide presentation for the class and the mentors. **Presented on 10/18.** Formal Project Proposal, in form of NIH proposal. | PS3 out on:L10-R8<br><br>**due Mon 10/21** | 6 | Tue, Oct 08 | Module III: Regulation, Epigenomics, Networks | Foundations ★ |
| | | | Thu, Oct 10 | | |
| | | | Fri, Oct 11 | | |
| | | 7 | Tue, Oct 15 | | Frontiers |
| | | | Thu, Oct 17 | | |
| | | | Fri, Oct 18 | | |
| | | | Thu, Oct 17 | | Project feedba ★ |
| Evaluate/discuss three peer proposals, NIH review format. **Reviews / Panel Discussion Mon 10/26. Written reviews due** | PS4 out on:L15-R10 | 8 | Tue, Oct 22 | Module IV: Evolution and Phylogenetics | Foundations |
| | | | Thu, Oct 24 | | ★ |
| | | | Fri, Oct 25 | | |
| | | | Fri, Oct 26 | | Panel Discuss ★ |
| Address peer evaluations, revise aims, scope, list of final deliverables / goals. Revised | | 9 | Tue, Oct 29 | | Frontiers |
| | **due Mon 11/04** | | Thu, Oct 31 | | |
| | | | Fri, Nov 01 | | |
| Continue making substantial progress on proposed milstones.Write outline of final report. **Midcourse progress report Due on Mon 11/18.** **Project final score projection from course staff by Friday** | PS5 out on:L15-R10<br><br>**due Mon 11/18** | 10 | Tue, Nov 05 | Module V: Population Genetics and Demography | Foundations |
| | | | Thu, Nov 07 | | |
| | | | Fri, Nov 08 | | |
| | | 11 | Tue, Nov 12 | | Frontiers |
| | | | Thu, Nov 14 | | |
| | | | Fri, Nov 15 | | |
| | | | Fri, Nov 15 | | Progress feed ★ |
| Complete your milestones, finalize results, figures, write-up in conference publication format. As part of report, comment on your overall project experience. **Final written report Due 12/7** | No more psets! (work on your final project) | 12 | Tue, Nov 19 | Quiz | Quiz |
| | | | Thu, Nov 21 | Module VI: Current Research Directions | Frontiers ★ |
| | | | Fri, Nov 22 | | |
| | | 13 | Tue, Nov 26 | | |
| | | | Thu, Nov 28 | | |
| | | | Fri, Nov 29 | | |
| | | | Mon, Dec 02 | | One-on-one m ★ |
| | | 14 | Tue, Dec 03 | | |
| | | | Thu, Dec 05 | | |
| | | | Fri, Dec 06 | Final Presentations | Frontiers ★ |
| Conference format slide presentation. **Talks on 12/10** | | 15 | Tue, Dec 10 | | ★ |
| | | | Tue, Dec 10 | | |

25

# Details on the final project

- Milestones ensure sufficient planning / feedback
  - Set-up: find project matching your skills and interests
  - Team: common interests and complementary skills
  - Inspiration: last year's projects, and recent papers
  - Proposal: establish milestones, deliverables, expectations
  - Midcourse: see endpoint, outline report, methods, figures
- Periodic mentoring sessions
  - Senior students and postdocs can serve as your mentors
  - Group discussions to share ideas, guidance, feedback
  - Peer-review: think critically about peer proposals, receive feedback/suggestions, respond to critiques, adjust course
- Real-world experience, condensed in a single term
  - Grant/fellowships proposals, peer review, yearly reports, budget time/effort, collaboration, paper writing, give talk

# Finding a research mentor / research advisor

- Chance to meet faculty at MIT/Broad/Harvard:
    - Through guest lectures and mentoring
    - Topics and papers covered in the lectures
    - Experts on: (1) human comparative genomics, (2) lincRNAs, (3) metabolic modeling, (4) disease mapping, selection, evolution and ecology (following four modules)
- Chance to meet senior students and postdocs:
    - On: coding genes, ncRNAs, regulatory motifs, networks, epigenomics, phylogenomics (again on each module)
    - Mentorship sessions with entire MIT CompBio group
- Your own personal research experience:
    - collaborators, datasets
    - learn active research directions, frontiers
    - living, breathing changing field

# Putting it all together

# Course Grading

- Grading:

| Problem sets 30% | Final Project 40% | Midterm 20% | Scrib10% |
|---|---|---|---|

- 4 problem sets:
  - Each problem set: 7-10%, covers 3-4 lectures, contains 3-4 problems.
  - Algorithmic problems and programming assignments (PS1 out now)
  - Graduate version includes additional problem on current research
- Final project
  - Introduction to research in computational biology (7 weeks!)
  - Includes peer-reviewed NIH-style proposal and much feedback
- Quiz
  - In-class quiz (Tue Nov 15). No final exam.
- Collaboration policy
  - Collaboration allowed, but you must:
    - Work independently on each problem before discussing it
    - Write solutions on your own
    - Acknowledge sources and collaborators. No outsourcing.

# Why <u>Computational</u> Biology ?

# Why <u>Computational</u> Biology:  Last year's answers

- Lots of data (* lots of data)
- There are rules
- Pattern finding
- It's *all* about data
- Ability to visualize
- Simulations, temporal relationships
- Guess + verify (generate hypotheses for testing)
- Propose mechanisms / theory to explain observations
- Networks / combinations of variables
- Efficiency (reduce experimental space to cover)
- Informatics infrastructure (ability to combine datasets)
- Correlations, higher-order relationships
- Cycle from hypothesis generation to testing condensed
- Life itself is digital.  Understand cellular instruction set

TATTGAATTTTCAAAAATTCTTACTTTTTTTTTGGATGGACGCAAAGAAGTTTAATAATCATATTACATGGCATTACCACCATATA
ATCCATATCTAATCTTACTTATATGTTGTGGAAATGTAAAGAGCCCCATTATCTTAGCCTAAAAAAACCTTCTCTTTGGAACTTT
AATACGCTTAACTGCTCATTGCTATATTGAAGTACGGATTAGAAGCCGCCGAGCGGGCGACAGCCCTCCGACGGAAGACTCTCCT
GCGTCCTCGTCTTCACCGGTCGCGTTCCTGAAACGCAGATGTGCCTCGCGCCGCACTGCTCCGAACAATAAAGATTCTACAATAC
TTTTATGGTTATGAAGAGGAAAAATTGGCAGTAACCTGGCCCCACAAACCTTCAAATTAACGAATCAAATTAACAACCATAGGAT
ATGCGATTAGTTTTTTAGCCTTATTTCTGGGGTAATTAATCAGCGAAGCGATGATTTTTGATCTATTAACAGATATATAAATGGA
CTGCATAACCACTTTAACTAATACTTTCAACATTTTCAGTTTGTATTACTTCTTATTCAAATGTCATAAAAGTATCAACAAAAAT
TAATATACCTCTATACTTTAACGTCAAGGAGAAAAAACTATAATGACTAAATCTCATTCAGAAGAAGTGATTGTACCTGAGTTCAA
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAGAAATTTATAAGCGCTTATGATGCTAAACCGG
TTGTTGCTAGATCGCCTGGTAGAGTCAATCTAATTGGTGAACATATTGATTATTGTGACTTCTCGGTTTTACCTTTAGCTATTGA
GATATGCTTTGCGCCGTCAAAGTTTTGAACGAGAAAAATCCATCCATTACCTTAATAAATGCTGATCCCAAATTTGCTCAAAGGAA
CGATTTGCCGTTGGACGGTTCTTATGTCACAATTGATCCTTCTGTGTCGGACTGGTCTAATTACTTTAAATGTGGTCTCCATGTTG
ACTCTTTTCTAAAGAAACTTGCACCGGAAAGGTTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTTCTGTGAGGGTGATGTACC
GGCAGTGGATTGTCTTCTTCGGCCGCATTCATTTGTGCCGTTGCTTTAGCTGTTGTTAAAGCGAATATGGGCCCTGGTTATCATA
CAAGCAAAATTTAATGCGTATTACGGTCGTTGCAGAACATTATGTTGGTGTTAACAATGGCGGTATGGATCAGGCTGCCTCTGTT
GTGAGGAAGATCATGCTCTATACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTTAAATTTCCGCAATTAAAAAACCATGA
AGCTTTGTTATTGCGAACACCCTTGTTGTATCTAACAAGTTTGAAACCGCCCCAACCAACTATAATTTAAGAGTGGTAGAAGTCA
AGCTGCAAATGTTTTAGCTGCCACGTACGGTGTTGTTTTACTTTCTGGAAAAGAAGGATCGAGCACGAATAAAGGTAATCTAAGA
TCATGAACGTTTATTATGCCAGATATCACAACATTTCCACACCCTGGAACGGCGATATTGAATCCGGCATCGAACGGTTAACAAA
CTAGTACTAGTTGAAGAGTCTCTCGCCAATAAGAAACAGGGCTTTAGTGTTGACGATGTCGCACAATCCTTGAATTGTTCTCGCG
ATTCACAAGAGACTACTTAACAACATCTCCAGTGAGATTTCAAGTCTTAAAGCTATATCAGAGGGCTAAGCATGTGTATTCTGAA
TAAGAGTCTTGAAGGCTGTGAAATTAATGACTACAGCGAGCTTTACTGCCGACGAAGACTTTTTCAAGCAATTTGGTGCCTTGAT
GAGTCTCAAGCTTCTTGCGATAAACTTTACGAATGTTCTTGTCCAGAGATTGACAAAATTTGTTCCATTGCTTTGTCAAATGGAT
TGGTTCCCGTTTGACCGGAGCTGGCTGGGGTGGTTGTACTGTTCACTTGGTTCCAGGGGGCCCAAATGGCAACATAGAAAAGGTA
AAGCCCTTGCCAATGAGTTCTACAAGGTCAAGTACCCTAAGATCACTGATGCTGAGCTAGAAAATGCTATCATCGTCTCTAAACCA
TTGGGCAGCTGTCTATATGAATTATAAGTATACTTCTTTTTTTTACTTTGTTCAGAACAACTTCTCATTTTTTTCTACTCATAACT
GCATCACAAAATACGCAATAATAACGAGTAGTAACACTTTTATAGTTCATACATGCTTCAACTACTTAATAAATGATTGTATGATA
TTTTCAATGTAAGAGATTTCGATTATCCACAAACTTTAAAACACAGGGACAAAATTCTTGATATGCTTTCAACCGCTGCGTTTTG
CCTATTCTTGACATGATATGACTACCATTTTGTTATTGTACGTGGGGCAGTTGACGTCTTATCATATGTCAAAGTCATTTGCGAA
TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAAGAGATTGCCGTCTTGAAACTTTTTGTCCTTTTTTTTTTCCGGGGACTCTAC
AACCCTTTGTCCTACTGATTAATTTTGTACTGAATTTGGACAATTCAGATTTTAGTAGACAAGCGCGAGGAGGAAAAGAAATGAC
AAATTCCGATGGACAAGAAGATAGGAAAAAAAAAAAGCTTTCACCGATTTCCTAGACCGGAAAAAAGTCGTATGACATCAGAATG
ATTTTCAAGTTAGACAAGGACAAAATCAGGACAAATTGTAAAGATATAATAAACTATTTGATTCAGCGCCAATTTGCCCTTTTCC
TCCATTAAATCTCTGTTCTCTCTTACTTATATGATGATTAGGTATCATCTGTATAAAACTCCTTTCTTAATTTCACTCTAAAGCA
CCATAGAGAAGATCTTTCGGTTCGAAGACATTCCTACGCATAATAAGAATAGGAGGGAATAATGCCAGACAATCTATCATTACAT
GCGGCTCTTCAAAAAGATTGAACTCTCGCCAACTTATGGAATCTTCCAATGAGACCTTTGCGCCAAATAATGTGGATTTGGAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTTATGAGGCATCGAGCTTTGAAGAAAAGTAAGCTCAGAAAAACCTCAAT
CTCATTCTGGAAGAAATCTATTATGAATATGTGGTCGTTGACAAATCAATCTTGGGTGTTTCTATTCTGGATTCATTTATGTAC
AGGACTTGAAGCCCGTCGAAAAGAAAGGCGGGTTTGGTCCTGGTACAATTATTGTTACTTCTGGCTTGCTGAATGTTTCAATAT
ACTTGGCAAATTGCAGCTACAGGTCTACAACTGGGTCTAAATTGGTGGCAGTGTTGGATAACAATTTGGATTGGGTACGGTTTCG

TATTGAATTTTCAAAAATTCTTACTTTTTTTTTTGGATGGACGCAAAGAAGTTTAATAATCATATTACATGGCATTACCACCATATA

ATCCATATCTAATCTTAC**TTATA**TGTTGTGGAAATGTAAAGAGCCCCATTATCTTAGCCTAAAAAAACCTTCTCTTTGGAACTTTC

AATACGCTTAACTGCTCATTGCTATATTGAAGTA**CGG**ATTAGAAGCCG**CCG**AG**CGG**GCGACAGCCCT**CCG**A**CGG**AAGACTCTCCTC

GCGTCCTCGTCTTCACCGGTCGCGTTCCTGAAACGCAGATGTGCCT**CGC**GCCGCACTGCT**CCG**AACAATAAAGATTCTACAATACT

TTTTATGGTTATGAAGAGGAAAAATTGGCAGTAACCTGG**CCCCA**CAAACCTTCAAATTAACGAATCAAATTAACAACCATAGGATG

ATGCGATTAGTTTTTTAGCCTTATTTC**TGGGG**TAATTAATCAGCGAAGCGATGATTTTTGATCTATTAACAGATA**TATAA**ATGGAA

CTGCATAACCACTTTAACTAATACTTTCAACATTTTCAGTTTGTATTACTTCTTATTCAAATGTCATAAAGTATCAACAAAAAT

TAATATACCTCTATACTTTAACGTCAAGGAGAAAAAACTATA**ATGACTAAATCTCATTCAGAAGAAGTGATTGTACCTGAGTTCAA**

**TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAGAAATTTATAAGCGCTTATGATGCTAAACCGG**

**TTGTTGCTAGATCGCCTGGTAGAGTCAATCTAATTGGTGAACATATTGATTATTGTGACTTCTCGGTTTTTACCTTTAGCTATTGAT**

**GATATGCTTTGCGCCGTCAAAGTTTTGAACGAGAAAAATCCATCCATTACCTTAATAAATGCTGATCCCAAATTTGCTCAAAGGAA**

**CGATTTGCCGTTGGACGGTTCTTATGTCACAATTGATCCTTCTGTGTCGGACTGGTCTAATTACTTTAAATGTGGTCTCCATGTTG**

**ACTCTTTTCTAAAGAAACTTGCACCGGAAAGGTTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTTCTGTGAGGGTGATGTACCA**

**GGCAGTGGATTGTCTTCTTCGGCCGCATTCATTTGTGCCGTTGCTTTAGCTGTTGTTAAAGCGAATATGGGCCCTGGTTATCATAT**

**CAAGCAAAATTTAATGCGTATTACGGTCGTTGCAGAACATTATGTTGGTGTTAACAATGGCGGTATGGATCAGGCTGCCTCTGTTT**

**GTGAGGAAGATCATGCTCTATACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTTAAATTTCCGCAATTAAAAAACCATGAA**

**AGCTTTGTTATTGCGAACACCCTTGTTGTATCTAACAAGTTTGAAACCGCCCCAACCAACTATAATTTAAGAGTGGTAGAAGTCAC**

**AGCTGCAAATGTTTTAGCTGCCACGTACGGTGTTGTTTTACTTTCTGGAAAAGAAGGATCGAGCACGAATAAAGGTAATCTAAGAG**

**TCATGAACGTTTATTATGCCAGATATCACAACATTTCCACACCCTGGAACGGCGATATTGAATCCGGCATCGAACGGTTAACAAAG**

**CTAGTACTAGTTGAAGAGTCTCTCGCCAATAAGAAACAGGGCTTTAGTGTTGACGATGTCGCACAATCCTTGAATTGTTCTCGCGA**

**ATTCACAAGAGACTACTTAACAACATCTCCAGTGAGATTTCAAGTCTTAAAGCTATATCAGAGGGCTAAGCATGTGTATTCTGAAT**

**TAAGAGTCTTGAAGGCTGTGAAATTAATGACTACAGCGAGCTTTACTGCCGACGAAGACTTTTTCAAGCAATTTGGTGCCTTGATG**

**GAGTCTCAAGCTTCTTGCGATAAACTTTACGAATGTTCTTGTCCAGAGATTGACAAAATTTGTTCCATTGCTTTGTCAAATGGATC**

**TGGTTCCCGTTTGACCGGAGCTGGCTGGGGTGGTTGTACTGTTCACTTGGTTCCAGGGGGCCCAAATGGCAACATAGAAAAGGTAA**

**AAGCCCTTGCCAATGAGTTCTACAAGGTCAAGTACCCTAAGATCACTGATGCTGAGCTAGAAAATGCTATCATCGTCTCTAAACCA**

**TTGGGCAGCTGTCTATATGAATTATAA**GTATACTTCTTTTTTTTTACTTTGTTCAGAACAACTTCTCATTTTTTTCTACTCATAACT

GCATCACAAAATACGCAATAATAACGAGTAGTAACACTTTTATAGTTCATACATGCTTCAACTACTTAATAAATGATTGTATGATA

TTTTCAATGTAAGAGATTTCGATTATCCACAAACTTTAAAACACAGGGACAAAATTCTTGATATGCTTTCAACCGCTGCGTTTTGC

CCTATTCTTGACATGATATGACTACCATTTTGTTATTGTACGTGGGGCAGTTGACGTCTTATCATATGTCAAAGTCATTTGCGAAG

TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAAGAGATTGCCGTCTTGAAACTTTTTGTCCTTTTTTTTTCCGGGGACTCTAC

AA**CCCTTTGT**CCTACTGATTAA**TTTTGTAC**TGAATTT**GGACAAT**TCAGATTTTAGTAGACAAGCGCGAGGAGGAAAAGAAATGACA

AAATTCCGATGGACAAGAAGATAGAAAAAAAAAAGCTTTCACCGATTTCCTAGACCGGAAAAAAGTCGTATGACATCAGAATGA

ATTTTCAAGTTAGA**CAAGGAC**AAAATCAGGACAAATTGTAAAGATATAATAAACTATTTGATTCAGCGCCAATTTGCCCTTTTTCCA

TCCATTAAATCTCTGTTCTCTCTTACTTATATGATGATTAGGTATCATCTG**TATAA**AACTCCTTTCTTAATTTCACTCTAAAGCAT

CCATAGAGAAGATCTTTCGGTTCGAAGACATTCCTACGCATAATAAGAATAGGAGGGAATA**ATGCCAGACAATCTATCATTACATT**

**GCGGCTCTTCAAAAAGATTGAACTCTCGCCAACTTATGGAATCTTCCAATGAGACCTTTGCGCCAAATAATGTGGATTTGGAAAAA**

**TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTTATGAGGCATCGAGCTTTGAAGAAAAAGTAAGCTCAGAAAAACCTCAATA**

**CTCATTCTGGAAGAAAATCTATTATGAATATGTGGTCGTTGACAAATCAATCTTGGGTGTTTCTATTCTGGATTCATTTATGTACA**

**AGGACTTGAAGCCCGTCGAAAAAGAAAGGCGGGTTTGGTCCTGGTACAATTATTGTTACTTCTGGCTTGCTGAATGTTTCAATATC**

**ACTTGGCAAATTGCAGCTACAGGTCTACAACTGGGTCTAAATTGGTGGCAGTGTTGGATAACAATTTGGATTGGGTACGGTTTCC**

# Extracting signal from noise

# The components of genomes and gene regulation



Enhancers | Promoters | Post-transcriptional control | microRNAs

**Goal: A systems-level understanding of genomes and gene regulation:**

- <u>The genome:</u> Map reads, align genes/genomes, assembly strategies
- <u>The genes:</u> Protein-coding exons, introns, non-coding RNA, RNA folding
- <u>The control regions:</u> Promoters, enhancers, insulators, chromatin states
- <u>The actual words:</u> Regulatory motifs, high-resolution accessibility maps
- <u>The regulators:</u> Transcription factors, chromatin modifiers, nucleosomes
- <u>The dynamics:</u> Changing maps between cell types, across development
- <u>The networks:</u> regulator→enhancer→target, ChIP-seq, correlated activity
- <u>The grammars:</u> TF/motif/mark combinations, predictive models
- <u>Human variation:</u> Human diversity, population genomics, linkage maps
- <u>Evolution:</u> Phylogenetics, phylogenomics, coalescent, human ancestry
- <u>GWAS/QTLs:</u> Genome variation ⇔ organismal/molecular phenotypes
- <u>Disease:</u> Personal (epi)genomics, pharmacogenomics, synthetic biology

**6.047/6.878/HST.507 - Fall 2015** - Lectures: Manolis Kellis.    All lectures on Tuesday/Thursday at 1pm-2:30m.
All homeworks due on Tuesday at 8pn

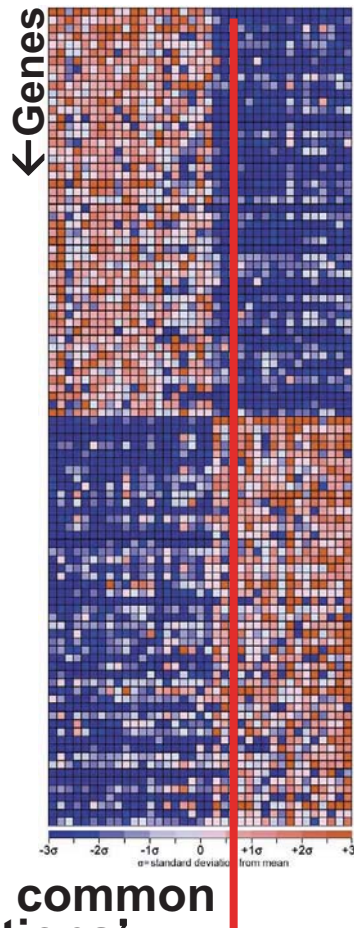| Project | Psets | Week | Date | Topic | | Lec | Topic | Read* |
|---|---|---|---|---|---|---|---|---|
| Describe your previous research, areas of interest in computational biology, type of project that best fits your interests. Post in a profile that lets your classmates know you and find potential partners. **Project profile due Tue 9/29** | PS1 out on:L1-L5 | 1 | Thu, Sep 10 | **Introduction** | | L1 | Intro: Biology, Algorithms, Machine Learning, Course Overview | 1 |
| | | | Fri, Sep 11 | | | R1 | Recitation 1: Biology and Probability Review | |
| | | 2 | Tue, Sep 15 | Module I: Aligning and Modeling Genomes | Foundations | L2 | Alignment I: Dynamic Programming, Global and local alignment | 2 |
| | | | Thu, Sep 17 | | | L3 | Alignment II: Database search, Rapid string matching, BLAST, BLOSUM | 3 |
| | | | Fri, Sep 18 | | | R2 | Recitation 2: Deriving Parameters of Alignment, Multiple Alignment | |
| | | 3 | Tue, Sep 22 | | Frontiers | L4 | Hidden Markov Models Part 1: Evaluation/Parsing, Viterbi, Forward algorithms | 7 |
| | due Tue 9/29 | | Thu, Sep 24 | | | L5 | Hidden Markov Models Part 2: Posterior Decoding, Learning, Baum-Welch | 8 |
| | | | Fri, Sep 25 | | | | No classes - student holiday | |
| | | | Fri, Sep 25 | | | Project Intro: about the projects, self introductions, mentor intro, example projects, teamwork 32D-507 | | |
| Identify previous project proposals, recent papers, and potential partners that match your areas of interest. List initial project ideas and partners. **Project area/team due Tue 10/6** | PS2 out on:L6-R4 | 4 | Tue, Sep 29 | Module II: Gene Expression and Networks | Foundations | L6 | Transcript structure: GenScan, RNA-seq, Mapping, De novo Assembly, Diff Expr | 12.3 |
| | | | Thu, Oct 1 | | | L7 | Expression Analysis: Clustering/Classification, K-means, Hierarchical, Bayesian | 15,16 |
| | | | Fri, Oct 2 | | | R3 | Recitation 3: Affinity Propagation Clustering and Random Forest Classification | |
| | due Tue 10/13 | 5 | Tue, Oct 6 | | Frontiers | L8 | Networks I: Bayesian inference, deep learning, network dynamics | 20,21 |
| Form teams of two, specify project goals, division of work, milestones, datasets, challenges Prepare slide presentation for the class and the mentors. **Project proposal due Tue 10/20. Presented on Fri 10/23** | | | Thu, Oct 8 | | | L9 | Networks II: Network learning, structure, spectral methods | 20,21 |
| | | | Fri, Oct 9 | | | R4 | Recitation 4: Small and Large Regulatory RNAs: lincRNA, miRNA, piRNA… | |
| | | | Fri, Oct 9 | | | Project Planning: research areas, initial ideas, type of project, mentor matching, finding partners 32D-507 | | |
| | PS3 out on:L10-R6 | 6 | Tue, Oct 13 | Module III: Gene Regulation & Epigenomics | | | No Classes - Monday Schedule - October 13, 2015 | |
| | | | Thu, Oct 15 | | Foundations | L10 | Regulatory Motifs: Discovery, Representation, PBMs, Gibbs Sampling, EM | 17 |
| | | | Fri, Oct 16 | | | R5 | Recitation 5: Gapped Motif Discovery, DNAShape, PBMs, Selex | |
| | due Tue 10/27 | 7 | Tue, Oct 20 | | Frontiers | L11 | Epigenomics: ChIP-Seq, Read mapping, Peak calling, IDR, Chromatin states | 19 |
| | | | Thu, Oct 22 | | | L12 | RNA modifications: RNA editing, Translation regulation, Splicing regulation | 11 |
| | | | Fri, Oct 23 | | | R6 | Recitation 6: Dimensionality Reduction | |
| | | | Fri, Oct 23 | | | Project feedback: Prepare 2-3 slide presentation of your term project for your mentor. 32D-507 at 4-5pm | | |
| Evaluate/discuss three peer proposals, NIH review format. **Review Panels Fri 10/30** **Reviews back Tue 11/3** | PS4 out on:L13-R8 | 8 | Tue, Oct 27 | Module IV: Population and Disease Genetics | Foundations | L13 | Resolving human ancestry and human history from genetic data | 29 |
| | | | Thu, Oct 29 | | | L14 | Disease Association Mapping, GWAS, organismal phenotypes | 31 |
| | | | Fri, Oct 30 | | | R7 | Recitation 7: Robinson-Foulds Distance and Coalescent Process | |
| | | | Fri, Oct 30 | | | Panel Discussion: reconciling critiques, strategies for improvement, feedback to author 32D-507 | | |
| Address peer evaluations, revise aims, scope, list of final deliverables / goals. **Response due Thu 11/12** | due Tue 11/10 | 9 | Tue, Nov 3 | | Frontiers | L15 | Quantitative trait mapping, molecular traits, eQTLs | 32 |
| | | | Thu, Nov 5 | | | L16 | Missing Heritability, Complex Traits, Interpret GWAS, Rank-based enrichment | 33 |
| | | | Fri, Nov 6 | | | R8 | Recitation 8: Suffix Trees and Arrays | |
| | PS5 out on:L17-R10 | 10 | Tue, Nov 10 | Module V: Comparative Genomics and Evolution | Foundations | | No lecture, veterans day holiday - Monday/Tuesday | |
| | | | Thu, Nov 12 | | | L17 | Comparative genomics and Evolutionary signatures | 4 |
| Continue making substantial progress on proposed milestones.Write outline of final report. **Midcourse report due Thu 11/19.** **Score projection 11/24** | | | Fri, Nov 13 | | | R9 | Recitation 9: Review of Phylogeny and Molecular Evolution | |
| | due Tue 12/1 | 11 | Tue, Nov 17 | | Frontiers | L18 | Phylogenetics: Molecular evolution, Tree building, Phylogenetic inference | 27 |
| | | | Thu, Nov 19 | | | L19 | Phylogenomics: Gene/species trees, reconciliation, recombination graphs | 28 |
| | | | Fri, Nov 20 | | | R10 | Recitation 10: Linkage Disequilibrium, Haplotype Phasing, and Genotype Imputation | |
| Complete your milestones, finalize results, figures, write-up in conference publication format. As part of report, comment on your overall project experience. **Written report due Sun 12/6** | No more psets! (work on your final project) | 12 | Tue, Nov 24 | Module VI: Current Research Directions | Frontiers | | In Class Quiz (the only quiz - the class has no final exam) - covers L1-R11 | |
| | | | Thu, Nov 26 | | | | No lecture, thanksgiving break - Thu Nov 26, 2015 | |
| | | | Fri, Nov 27 | | | | No recitation, thanksgiving break | |
| | | 13 | Tue, Dec 1 | | | L20 | Personal Genomics, Disease Epigenomics: Systems approaches to disease | 34,36 |
| | | | Thu, Dec 3 | | | L21 | Three-dimensional chromatin interactions: 3C, 5C, HiC, ChIA-Pet | 30 |
| | | | Fri, Dec 4 | | | R11 | Recitation 11: Project Tips - Write-up, Slides, Final Presentation in 32D-507 | |
| Conference format slide pres. **Talks on Thu 12/10** | | 14 | Tue, Dec 8 | | | L22 | Genome Engineering with CRISPR/Cas9 and related technologies | |
| | | | Thu, Dec 10 | | | | Final Presentations - Part I (1pm). 32-141 | |
| | | | Thu, Dec 10 | | | | Final Presentations - Part II (3pm). 32D-507 | |

\* readings refer to chapters in compiled 2014 scribe notes, available in the materials folder
\*\* recitation topics will be adjusted to respond to lecture and student needs

# Coupling each topic with foundational CS tools

| Lect | Fundamental bio problem | Foundational comp. tool |
|---|---|---|
| 1 | Introduction | |
| 2 | Sequence alignment | Dynamic programming |
| 3 | Database search | Hashing |
| 4,5 | Modeling biological signals | HMMs/Modeling/Learning/EM |
| 6,7 | Transcriptome analysis | Clustering / EM |
| 8,9 | Regulatory networks | Graph algorithms, spectral analysis |
| 10 | Regulatory motifs | Information/Gibbs Sampling/EM |
| 11 | Epigenomics | Classification / Modeling |
| 13-16 | Population Genetics | Statistical modeling and inference |
| 18-19 | Gene trees and species trees | Phylogenetics/Bayesian inference |

# Overview of the 5 modules

# Challenges in Computational Biology

(4) Genome Assembly

(5) Regulatory motif discovery

(1) Gene Finding

DNA

(2) Sequence alignment

(6) Comparative Genomics

(7) Evolutionary Theory

```
TCATGCTAT
TCGTGATAA
TGAGGATAT
TTATCATAT
TTATGATTT
```

(3) Database lookup

(8) Gene expression analysis

RNA transcript

(9) Cluster discovery

(10) Gibbs sampling

(11) Protein network analysis

(12) Metabolic modelling

(13) Emerging network properties

# Module 1: Aligning and Modeling Genomes

| 1 | Thu, Sep 10 | **Introduction** | | L1 | Intro: Biology, Algorithms, Machine Learning, Course Overview | 1 |
| | Fri, Sep 11 | | | R1 | Recitation 1: Biology and Probability Review | |
| 2 | Tue, Sep 15 | Module I: Aligning and Modeling Genomes | Foundations | L2 | Alignment I: Dynamic Programming, Global and local alignment | 2 |
| | Thu, Sep 17 | | | L3 | Alignment II: Database search, Rapid string matching, BLAST, BLOSUM | 3 |
| | Fri, Sep 18 | | | R2 | Recitation 2: Deriving Parameters of Alignment, Multiple Alignment | |
| 3 | Tue, Sep 22 | | Frontiers | L4 | Hidden Markov Models Part 1: Evaluation/Parsing, Viterbi, Forward algorithms | 7 |
| | Thu, Sep 24 | | | L5 | Hidden Markov Models Part 2: Posterior Decoding, Learning, Baum-Welch | 8 |
| | Fri, Sep 25 | | | | No classes - student holiday | |
| | Fri, Sep 25 | | | | Project Intro: about the projects, self introductions, mentor intro, example projects, teamwork 32D-507 | |

- Foundations vs. frontiers
  - Foundations: Classical computational methods / biological topics
  - Frontiers: Latest developments, open questions, research areas
  - Duality for each:  basic problems / fundamental techniques

- Sequence alignment:
  - Local/global alignment: infer nucleotide-level evolutionary events
  - Database search: scan for regions that may have common ancestry

- Hidden Markov Models
  - Hidden Markov Models (HMMs): Central tool in CS
  - Decoding, evaluation, parsing, likelihood, scoring

# Dynamic Programming Algorithms: Align, HMMs



- ## Sequence alignment
- ## Hidden Markov Models
- ## DP: Core computational technique
  - Pervasive in computer science, and computational biology
  - Fully explore exponential search spaces in poly time!
  - Greedy algorithms will not work, back-tracking, saving soln
  - Special requirements: Optimal substructure
  - Found in: alignment, HMMs, phylogeny, genetics, pop gen…

# Module II: Gene expression analysis and transcripts

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| PS2 out on:L6-R4 | 4 | Tue, Sep 29 | Module II: Gene Expression and Networks | Foundations | L6 | Transcript structure: GenScan, RNA-seq, Mapping, De novo Assembly, Diff Expr | 12.3 |
| | | Thu, Oct 1 | | | L7 | Expression Analysis: Clustering/Classification, K-means, Hierarchical, Bayesian | 15,16 |
| | | Fri, Oct 2 | | | R3 | Recitation 3: Affinity Propagation Clustering and Random Forest Classification | |
| | 5 | Tue, Oct 6 | | Frontiers | L8 | Networks I: Bayesian inference, deep learning, network dynamics | 20,21 |
| | | Thu, Oct 8 | | | L9 | Networks II: Network learning, structure, spectral methods | 20,21 |
| due | | Fri, Oct 9 | | | R4 | Recitation 4: Small and Large Regulatory RNAs: lincRNA, miRNA, piRNA… | |
| Tue 10/13 | | Fri, Oct 9 | | Project Planning: research areas, initial ideas, type of project, mentor matching, finding partners 32D-507 | | | |

- ## Computational foundations:
  – Unsupervised Learning: Expectation Maximization
  – Supervised learning: generative/discriminative models
  – Read mapping, significance testing, splice graphs
- ## Biological frontiers:
  – PS2: Modeling conservation, GC content, CpG islands
  – L6/L7: Genome annotation and parsing
  – L8: Gene expression analysis: cluster genes/conditions
  – L9: Regulatory motif discovery: EM, gibbs sampling, info

# Natural 1st step: group similar rows/columns
# Clustering

➜ **Similar cell types**

Conditions➜

← Genes

Armstrong, Nature Gen 2002

**Reveal common 'conditions'**

➜ **Similarly-behaving groups of genes**

Conditions➜

← Genes

Alizadeh, Nature 2000

**Reveal common gene behaviors**

# If labels are known: find more of same type
# Classification

➔ Classify diseases      ➔ Classify genes in different pathways



Armstrong, Nature Gen 2002

Alizadeh, Nature 2000

**Find features that distinguish known classes**

**Find additional members of existing gene classes**
**Predict function of uncharacterized genes**

45

# Module III: Epigenomics and gene regulation

| | | | | | | |
|---|---|---|---|---|---|---|
| **6** | Tue, Oct 13 | Module III: Gene Regulation & Epigenomics | Foundations | | No Classes - Monday Schedule - October 13, 2015 | |
| | Thu, Oct 15 | | | L10 | Regulatory Motifs: Discovery, Representation, PBMs, Gibbs Sampling, EM | 17 |
| | Fri, Oct 16 | | | R5 | Recitation 5: Gapped Motif Discovery, DNAShape, PBMs, Selex | |
| **7** | Tue, Oct 20 | | Frontiers | L11 | Epigenomics: ChIP-Seq, Read mapping, Peak calling, IDR, Chromatin states | 19 |
| | Thu, Oct 22 | | | L12 | RNA modifications: RNA editing, Translation regulation, Splicing regulation | 11 |
| | Fri, Oct 23 | | | R6 | Recitation 6: Dimensionality Reduction | |
| | Fri, Oct 23 | | | | Project feedback: Prepare 2-3 slide presentation of your term project for your mentor. 32D-507 at 4-5pm | |

- Computational Foundations
  - Hidden Markov Models (HMMs): Central tool in CS
  - Decoding, evaluation, parsing, likelihood, scoring
  - Unsupervised Learning: Expectation Maximization
  - Supervised learning: generative/discriminative models
- Biological frontiers:
  - PS2: Modeling conservation, GC content, CpG islands
  - L6/L7: Genome annotation and parsing
  - L8: Gene expression analysis: cluster genes/conditions
  - L9: Regulatory motif discovery: EM, gibbs sampling, info

# Motifs summarize TF sequence specificity

| Target genes bound by ABF1 regulator | | Coordinates | | Genome sequence at bound site |
|---|---|---|---|---|
| ACS1 | acetyl CoA synthetase | -491 | -479 | \|ATCATTCTGGACG\| |
| ACS1 | acetyl CoA synthetase | -433 | -421 | \|ATCATCTCGGACG\| |
| ACS1 | acetyl CoA synthetase | -311 | -299 | \|ATCATTTGCCACG\| |
| CHA1 | catabolic L-serine dehydratase | -280 | -254 | A\|ATCACCGCGAACG\|GA |
| ENO2 | Enolase | -470 | -461 | ggcgttat\|GTCACTAACGACG\|tgcacca |
| HMR | silencer | -256 | -283 | ATCAATAC\|ATCATAAAATACG\|AACGATC |
| LPD1 | lipoamide dehydrogenase | -288 | -300 | gat\|ATCAAAATTAACG\|tag |
| LPD1 | lipoamide dehydrogenase | -301 | -313 | gat\|ATCACCGTTGACG\|tca |
| PGK | phosphoglycerate kinase | -523 | -496 | CAAACAA\|ATCACGAGCGACG\|GTAATTTC |
| RPC160 | RNA pol III/C 160 kDa subunit | -385 | -349 | \|ATCACTATATACG\|TGAA |
| RPC40 | RNA pol III/C 40 kDa subunit | -137 | -116 | \|GTCACTATAAACG\| |
| rpL2 | ribosomal protein L2 | -185 | -167 | TAAT\|aTCAcgtcACACG\|AC |
| SPR3 | CDC3/10/11/12 family homolog | -315 | -303 | \|ATCACTAAATACG\| |
| YPT1 | TUB2 | -193 | -172 | CCTAG\|GTCACTGTACACG\|TATA |

| Position | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Position Weight Matrix (PWM) | A | 56 | 4 | 4 | 81 | 4 | 23 | 15 | 27 | 31 | 31 | 89 | 23 | 4 | 58 |
| | G | 32 | 4 | 4 | 12 | 4 | 31 | 23 | 4 | 19 | 23 | 4 | 4 | 89 | 35 |
| | C | 4 | 4 | 89 | 4 | 58 | 12 | 23 | 19 | 19 | 23 | 4 | 69 | 4 | 4 |
| | T | 4 | 89 | 4 | 4 | 35 | 35 | 39 | 50 | 31 | 23 | 4 | 4 | 4 | 4 |
| Motif Logo | | | | | | | | | | | | | | | |
| Consensus | | R | T | C | A | Y | N | N | H | N | N | A | C | G | R |

- Summarize information

- Integrate many positions

- Measure of information

- Distinguish motif vs. motif instance

- Assumptions:
  - Independence
  - Fixed spacing

47

# Starting positions ⇔ Motif matrix

- given <u>aligned</u> sequences ➡ easy to compute profile matrix

**shared motif**

**sequence positions**

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| **A** | 0.1 | 0.3 | 0.1 | 0.2 | 0.2 | 0.4 | 0.3 | 0.1 |
| **C** | 0.5 | 0.2 | 0.1 | 0.1 | 0.6 | 0.1 | 0.2 | 0.7 |
| **G** | 0.2 | 0.2 | 0.6 | 0.5 | 0.1 | 0.2 | 0.2 | 0.1 |
| **T** | 0.2 | 0.3 | 0.2 | 0.2 | 0.1 | 0.3 | 0.3 | 0.1 |

**given profile matrix**

- **easy to find starting position probabilities**

**Key idea:  Iterative procedure for estimating both, given uncertainty**

**(learning problem with hidden variables:  the starting positions)**

# Multivariate HMM for Chromatin States



Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Ernst, Jason and Manolis Kellis. "Discovery and characterization of chromatin states for systematic annotation of the human genome." Nature Biotechnology 28, no. 8 (2010): 817-825.

**Ernst and Kellis**
**Nature Biotech 2010**

# Modules IV and V: Evolution/phylogeny/populations

| | | | | | | |
|---|---|---|---|---|---|---|
| 8 | Tue, Oct 27 | Module IV: Population and Disease Genetics | Foundations | L13 | Resolving human ancestry and human history from genetic data | 29 |
| | Thu, Oct 29 | | | L14 | Disease Association Mapping, GWAS, organismal phenotypes | 31 |
| | Fri, Oct 30 | | | R7 | Recitation 7: Robinson-Foulds Distance and Coalescent Process | |
| | Fri, Oct 30 | | Panel Discussion: reconciling critiques, strategies for improvement, feedback to author 32D-507 | | | |
| 9 | Tue, Nov 3 | | Frontiers | L15 | Quantitative trait mapping, molecular traits, eQTLs | 32 |
| | Thu, Nov 5 | | | L16 | Missing Heritability, Complex Traits, Interpret GWAS, Rank-based enrichment | 33 |
| | Fri, Nov 6 | | | R8 | Recitation 8: Suffix Trees and Arrays | |
| 10 | Tue, Nov 10 | Module V: Comparative Genomics and Evolution | Foundations | | No lecture, veterans day holiday - Monday/Tuesday | |
| | Thu, Nov 12 | | | L17 | Comparative genomics and Evolutionary signatures | 4 |
| | Fri, Nov 13 | | | R9 | Recitation 9: Review of Phylogeny and Molecular Evolution | |
| 11 | Tue, Nov 17 | | Frontiers | L18 | Phylogenetics: Molecular evolution, Tree building, Phylogenetic inference | 27 |
| | Thu, Nov 19 | | | L19 | Phylogenomics: Gene/species trees, reconciliation, recombination graphs | 28 |
| | Fri, Nov 20 | | | R10 | Recitation 10: Linkage Disequilibrium, Haplotype Phasing, and Genotype Imputation | |

- ## Phylogenetics / Phylogenomics
  - Phylogenetics: Evolutionary models, Tree building, Phylo inference
  - Phylogenomics: gene/species trees, reconciliation, coalescent, pops

- ## Population genomics:
  - Learning population history from genetic data (David Reich)
  - Statistical genetics: disease mapping in populations (Mark Daly)
  - Measuring natural selection in human populations (Pardis Sabeti)
  - The missing heritability in genome-wide associations (Yaniv Erlich)

- ## And we're done! Last pset Nov 21st, In-class quiz on Nov 22nd
  - No lab 4! Then entire focus shifts to projects, Thanksgiving, Frontiers

# Characterizing sub-threshold variants in heart arrhythmia

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Arking, Dan E. et al. "Genetic association study of QT interval highlights role for calcium signaling pathways in myocardial repolarization." Nature Genetics 46, no. 8 (2014): 826-836.

*Focus on sub-threshold variants*
*(e.g. rs1743292 P=$10^{-4.2}$)*

*Trait: QRS/QT interval*

(1) Large cohorts, (2) many known hits
(3) well-characterized tissue drivers

# Evidence of Neanderthal→Human gene flow



NO
GENE FLOW

GENE FLOW

Neand.-African tMRCA
~825 kya

Average human tMRCA
~ 500 kya

Population split
~ 350 kya

~ 100 kya

Neandertal    Human1    Human2

Neandertal    Human1    Human2

hg18
BAC

hg18
BAC

Courtesy of Luna04 on wikipedia.
License: CC BY.

Courtesy of Luna04 on wikipedia.
License: CC BY.

Human-human divergence is

AVERAGE

Human-human divergence is

HIGH

# Structure of genetic code ⇔ evolutionary signatures

- Substitutions that preserve AA properties tolerated in coding exons
- Leads to specific evolutionary signatures associated with protein-coding genes
- The code itself could be rediscovered simply based on observed substitution patterns



$Q_C$ estimated from known coding regions



$Q_N$ estimated from non-coding regions

These specify different rates of codon substitution, which in turn lead to different probabilities of any given alignment:





$$Pr(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \frac{1}{10^{117}}$$

$$Pr(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \frac{1}{10^{152}}$$

$$Pr(\text{Leaves}; \mathbf{Q}_C, \underline{t}) = \frac{1}{10^{275}}$$

$$Pr(\text{Leaves}; \mathbf{Q}_N, \underline{t}) = \frac{1}{10^{254}}$$

# Distance matrix ⇔ Phylogenetic tree

|       | Hum | Mou | Rat | Dog | Cat |
|-------|------|-------|-------|-------|-------|
| **Human** | 0 | 4 | 5 | 7 | 6 |
| **Mouse** | h.y.m | 0 | 3 | 8 | 5 |
| **Rat** | h.y.r | m.r | 0 | 9 | 7 |
| **Dog** | h.z.x.d | m.y.z.x.d | r.y.z.x.d | 0 | 2 |
| **Cat** | h.z.x.c | m.y.z.x.c | r.y.z.x.c | d.c | 0 |

Tree implies a distance matrix $M_{ij}$

Map distances $D_{ij}$ to a tree

$$\min \sum_{ij} (D_{ij} - M_{ij})^2$$

Goal:

Minimize discrepancy between observed distances and tree-based distances

54

# 'Peeling' algorithm for P(D|B,T) term

$x_9 = $ "AAACTG"

$$P(x_1, ..., x_{2n-1}|T, t) = P(x_1|x_2, ..., x_{2n-1}, T, t)P(x_2|x_3, ..., x_{2n-1}, T, t)...P(x_{2n-1}|T, t)$$
$$= P(x_1|x_{parent(1)}, t_1)P(x_2|x_{parent(2)}, t_2)...P(x_{2n-1})$$
$$= P(x_{2n-1}) \prod_{i=1}^{2n-2} P(x_i|x_{parent(i)}, t_i)$$

1. Assume **sites j evolve independently**.
   - ➔ Treat each column of the alignment in isolation

2. Assume **branch independence**, conditioned on parent
   - ➔ Expand total joint probability into prod of $P(x_i|x_{parent}, t_i)$
   - ➔ Only $P(x_{2n-1})$ remains, root prior, background nucl. freq.

3. We know how to compute **$P(x_i|x_{parent(i)}, t_i)$** for fixed pair
   - ➔ Defined by our sequence model (JC, K2P, HKY, etc)
   - ➔ Easily calculate for any given assignment of internal nodes

4. As internal node values are not known ➔ **marginalize**
   - ➔ Sum over all possible values of all internal/root nodes
   - ➔ Let $x_{n+1}, ..., x_{2n-1}$ represent seqs of n-1 internal nodes

# Two types of gene-tree species-tree reconciliation

**Gene tree**

**Coalescence**

**Duplication & Loss**

➔ **DLCoal**

Species tree

A    B    C

A    B    C

A    B    C

A    B    C

- **Coalescent models of alleles in populations**
  - **Deal with 1-to-1 orthologs**
  - **Estimate divergence times, pop sizes, etc**
  - **Models move backward in time**
  - **Cannot cope with duplication and loss**

- **DL models of genes in species**
  - **Deal with paralogous families**
  - **Estimate birth death rates**
  - **Models move forward in time**
  - **Cannot cope with incomplete lineage sorting**

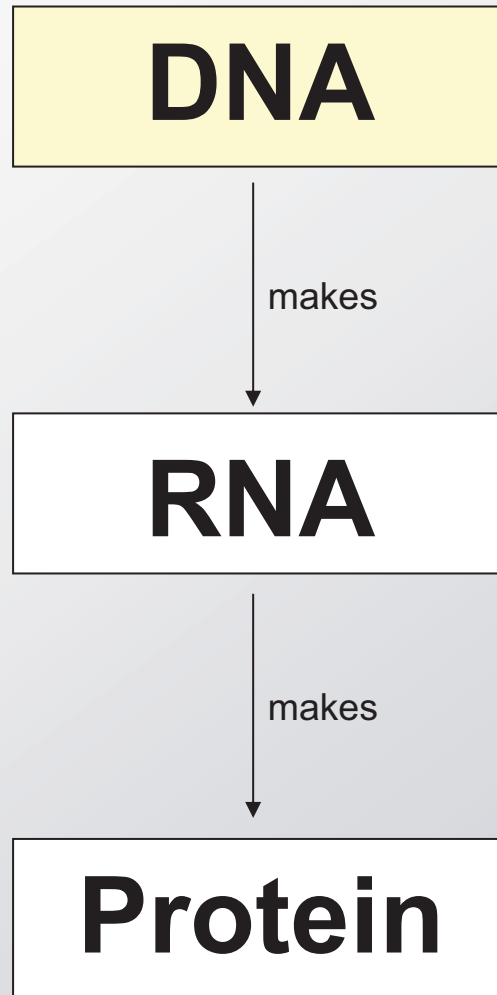| Project | Psets | Week | Date | Topic | | Lec | Topic |
|---|---|---|---|---|---|---|---|
| Describe your previous research, areas of interest in computational biology, type of project that best fits your interests. Post in a profile that lets your classmates know you and find potential partners. **Project profile due Tue 9/29** | PS1 out on:L1-L5 | 1 | Thu, Sep 10 | **Introduction** | | L1 | Intro: Biology, Algorithms, Machine Learning, Course Overview |
| | | | Fri, Sep 11 | | | R1 | Recitation 1: Biology and Probability Review |
| | | 2 | Tue, Sep 15 | Module I: Aligning and Modeling Genomes | Foundations | L2 | Alignment I: Dynamic Programming, Global and local alignment |
| | | | Thu, Sep 17 | | | L3 | Alignment II: Database search, Rapid string matching, BLAST, BLOSUM |
| | | | Fri, Sep 18 | | | R2 | Recitation 2: Deriving Parameters of Alignment, Multiple Alignment |
| | | 3 | Tue, Sep 22 | | Frontiers | L4 | Hidden Markov Models Part 1: Evaluation/Parsing, Viterbi, Forward algorithms |
| | | | Thu, Sep 24 | | | L5 | Hidden Markov Models Part 2: Posterior Decoding, Learning, Baum-Welch |
| | due Tue 9/29 | | Fri, Sep 25 | | | | No classes - student holiday |
| | | | Fri, Sep 25 | | | | Project Intro: about the projects, self introductions, mentor intro, example projects, teamwork 32D-50 |
| Identify previous project proposals, recent papers, and potential partners that match your areas of interest. List initial project ideas and partners. **Project area/team due Tue 10/6** | PS2 out on:L6-R4 | 4 | Tue, Sep 29 | Module II: Gene Expression and Networks | Foundations | L6 | Transcript structure: GenScan, RNA-seq, Mapping, De novo Assembly, Diff Expr |
| | | | Thu, Oct 1 | | | L7 | Expression Analysis: Clustering/Classification, K-means, Hierarchical, Bayesian |
| | | | Fri, Oct 2 | | | R3 | Recitation 3: Affinity Propagation Clustering and Random Forest Classification |
| | | 5 | Tue, Oct 6 | | Frontiers | L8 | Networks I: Bayesian inference, deep learning, network dynamics |
| | | | Thu, Oct 8 | | | L9 | Networks II: Network learning, structure, spectral methods |
| Form teams of two, specify project goals, division of work, milestones, datasets, challenges Prepare slide presentation for the class and the mentors. **Project proposal due Tue 10/20. Presented on Fri 10/23** | due Tue 10/13 | | Fri, Oct 9 | | | R4 | Recitation 4: Small and Large Regulatory RNAs: lincRNA, miRNA, piRNA... |
| | | | Fri, Oct 9 | | | | Project Planning: research areas, initial ideas, type of project, mentor matching, finding partners 32D- |
| | PS3 out on:L10-R6 | 6 | Tue, Oct 13 | Module III: Gene Regulation & Epigenomics | | | No Classes - Monday Schedule - October 13, 2015 |
| | | | Thu, Oct 15 | | Foundations | L10 | Regulatory Motifs: Discovery, Representation, PBMs, Gibbs Sampling, EM |
| | | | Fri, Oct 16 | | | R5 | Recitation 5: Gapped Motif Discovery, DNAShape, PBMs, Selex |
| | | 7 | Tue, Oct 20 | | Frontiers | L11 | Epigenomics: ChIP-Seq, Read mapping, Peak calling, IDR, Chromatin states |
| | | | Thu, Oct 22 | | | L12 | RNA modifications: RNA editing, Translation regulation, Splicing regulation |
| | due Tue 10/27 | | Fri, Oct 23 | | | R6 | Recitation 6: Dimensionality Reduction |
| | | | Fri, Oct 23 | | | | Project feedback: Prepare 2-3 slide presentation of your term project for your mentor. 32D-507 at 4-5 |
| Evaluate/discuss three peer proposals, NIH review format. **Review Panels Fri 10/30 Reviews back Tue 11/3** | PS4 out on:L13-R8 | 8 | Tue, Oct 27 | Module IV: Population and Disease Genetics | Foundations | L13 | Resolving human ancestry and human history from genetic data |
| | | | Thu, Oct 29 | | | L14 | Disease Association Mapping, GWAS, organismal phenotypes |
| | | | Fri, Oct 30 | | | R7 | Recitation 7: Robinson-Foulds Distance and Coalescent Process |
| | | | Fri, Oct 30 | | | | Panel Discussion: reconciling critiques, strategies for improvement, feedback to author 32D-507 |
| Address peer evaluations, revise aims, scope, list of final deliverables / goals. **Response due Thu 11/12** | due Tue 11/10 | 9 | Tue, Nov 3 | | Frontiers | L15 | Quantitative trait mapping, molecular traits, eQTLs |
| | | | Thu, Nov 5 | | | L16 | Missing Heritability, Complex Traits, Interpret GWAS, Rank-based enrichment |
| | | | Fri, Nov 6 | | | R8 | Recitation 8: Suffix Trees and Arrays |
| | PS5 out on:L17-R10 | 10 | Tue, Nov 10 | Module V: Comparative Genomics and Evolution | | | No lecture, veterans day holiday - Monday/Tuesday |
| | | | Thu, Nov 12 | | Foundations | L17 | Comparative genomics and Evolutionary signatures |
| Continue making substantial progress on proposed milstones. Write outline of final report. **Midcourse report due Thu 11/19. Score projection 11/24** | | | Fri, Nov 13 | | | R9 | Recitation 9: Review of Phylogeny and Molecular Evolution |
| | | 11 | Tue, Nov 17 | | Frontiers | L18 | Phylogenetics: Molecular evolution, Tree building, Phylogenetic inference |
| | due Tue 12/1 | | Thu, Nov 19 | | | L19 | Phylogenomics: Gene/species trees, reconciliation, recombination graphs |
| | | | Fri, Nov 20 | | | R10 | Recitation 10: Linkage Disequilibrium, Haplotype Phasing, and Genotype Imputatio |
| Complete your milestones, finalize results, figures, write-up in conference publication format. As par of report, comment on your overall project experience. **Written report due Sun 12/6** | No more psets! (work on your final project) | 12 | Tue, Nov 24 | Module VI: Current Research Directions | | | In Class Quiz (the only quiz - the class has no final exam) - covers L1-R11 |
| | | | Thu, Nov 26 | | | | No lecture, thanksgiving break - Thu Nov 26, 2015 |
| | | | Fri, Nov 27 | | | | No recitation, thanksgiving break |
| | | 13 | Tue, Dec 1 | | Frontiers | L20 | Personal Genomics, Disease Epigenomics: Systems approaches to disease |
| | | | Thu, Dec 3 | | | L21 | Three-dimensional chromatin interactions: 3C, 5C, HiC, ChIA-Pet |
| | | | Fri, Dec 4 | | | R11 | Recitation 11: Project Tips - Write-up, Slides, Final Presentation in 32D-507 |
| | | 14 | Tue, Dec 8 | | | L22 | Genome Engineering with CRISPR/Cas9 and related technologies |
| Conference format slide pres. **Talks on Thu 12/10** | | | Thu, Dec 10 | | | | Final Presentations - Part I (1pm). 32-141 |
| | | | Thu, Dec 10 | | | | Final Presentations - Part II (3pm). 32D-507 |

\* readings refer to chapters in compiled 2014 scribe notes, available in the materials folder
\*\* recitation topics will be adjusted to respond to lecture and student needs
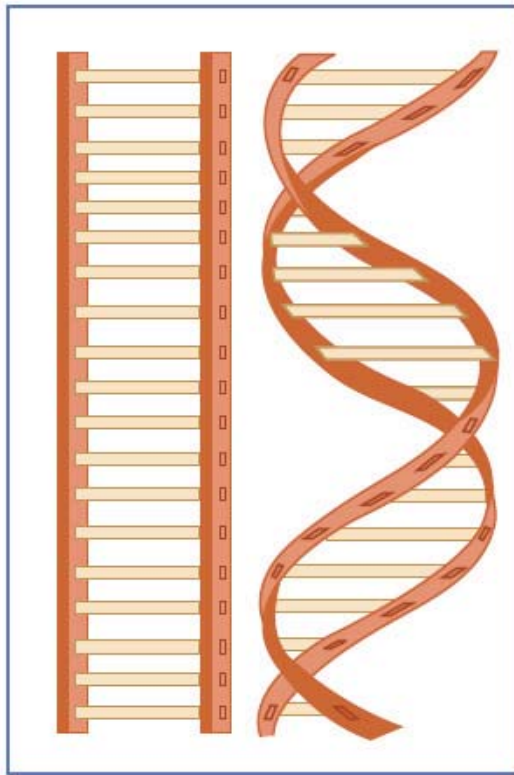
# Biology primer

Quick introduction to molecular biology
and information transfer within the cell

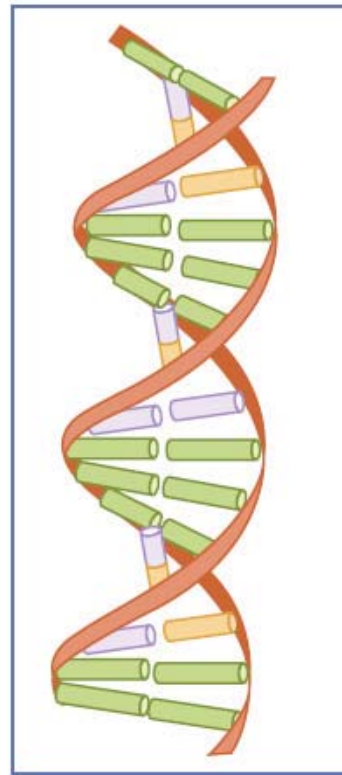# "Central dogma" of Molecular Biology

# DNA:  The double helix
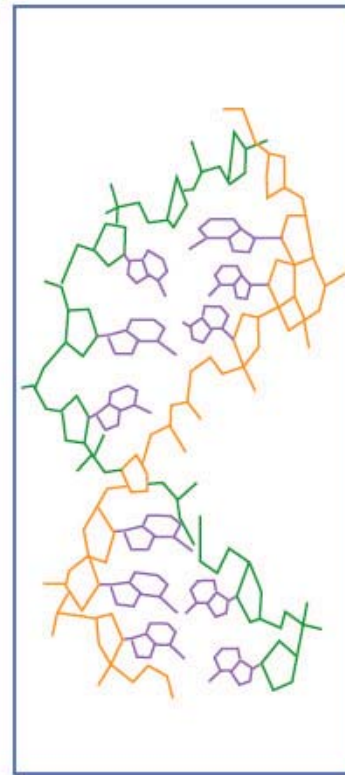
- The most noble molecule of our time



In fact, the two DNA strands are twisted around each other to make a double helix.
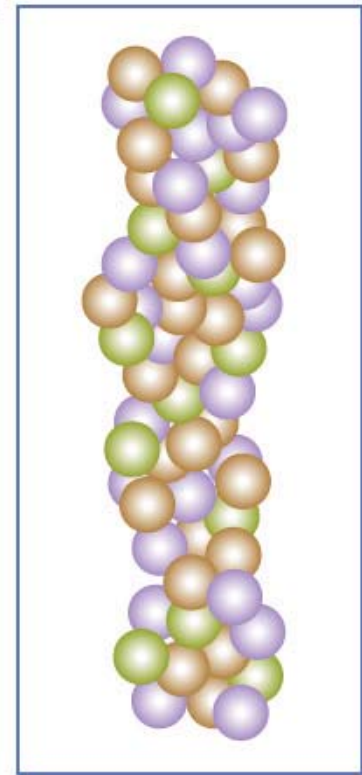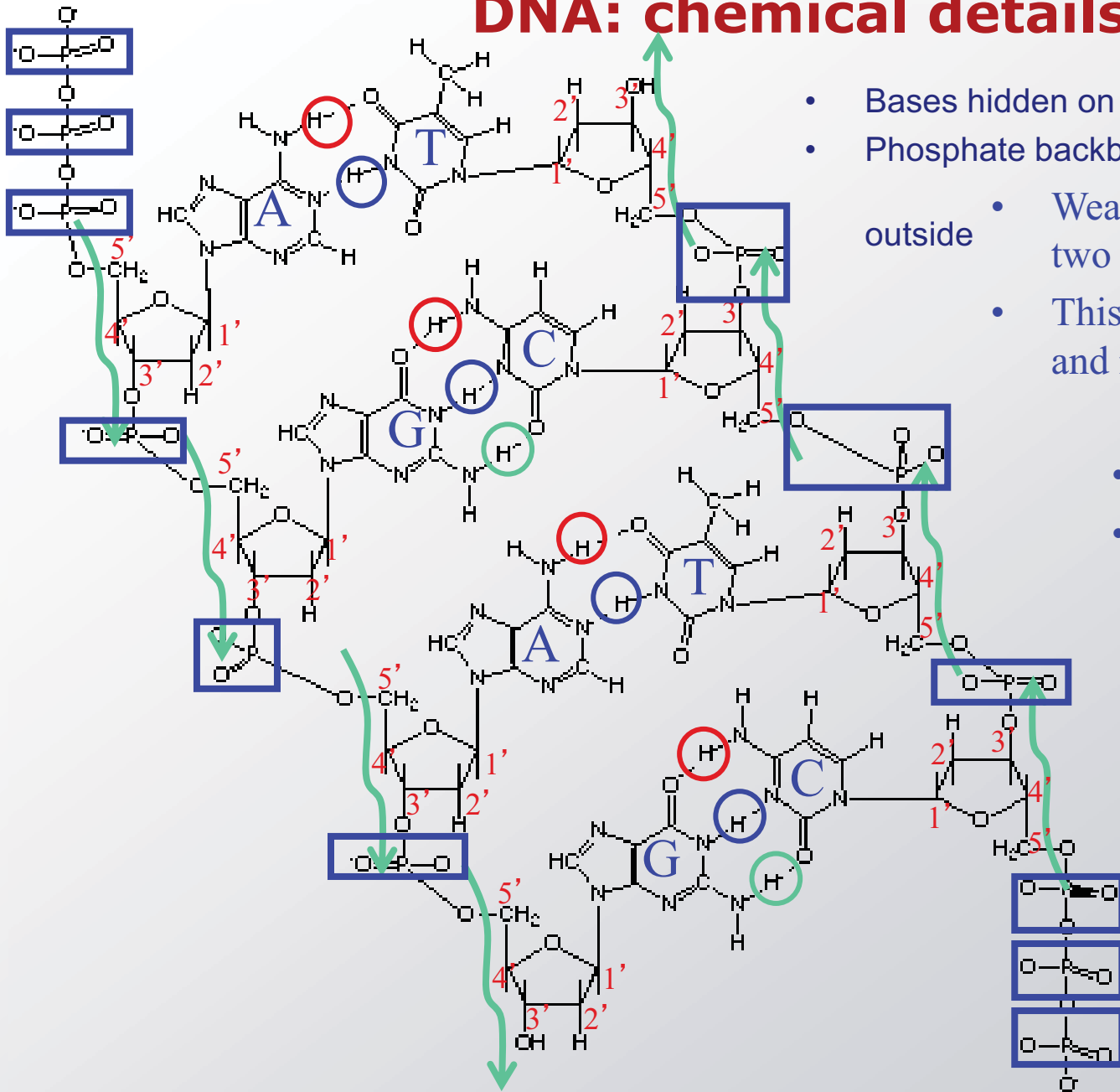
Traditional

Fancy

Chemical

Atomic

Image by MIT OpenCourseWare.

# DNA: the molecule of heredity

- Self-complementarity sets molecular basis of heredity
    - Knowing one strand, creates a template for the other
    - "It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material." Watson & Crick, 1953



Image by MIT OpenCourseWare.

# DNA: chemical details

- Bases hidden on the inside
- Phosphate backbone outside

- Weak hydrogen bonds hold the two strands together
- This allows low-energy opening and re-closing of two strands

- Anti-parallel strands
- Extension 5' → 3' tri-phosphate coming from newly added nucleotide

The only parings are:
- A with T
- C with G

# DNA: the four bases

## The Nucleotides of DNA



| Adenine | Guanosine | Thymine | Cytosine |
|---------|-----------|---------|----------|
| Purine | Purine | | |
| | | Pyrimidine | Pyrimidine |
| Weak | | Weak | |
| | Strong | | Strong |
| Amino | | | Amino |
| | Keto | Keto | |

# Alignment: all species/genes share common ancestry

Slide credit: Serafim Batzoglou

# Tree of Life



Cellular organisms without cell nuclei are Prokaryotes ("before kernel")

First Life-form

Chromalveolates

Rhizaria

Green Plants

Bryophytes

Embryophytes

Cell Nucleus and Mitochondria (Eukarya) "good kernel"

Vascular Plants (xylem)

Plants

Seeds (Spermatophytes)

Eudicots

Dicots

Enclosed Seeds Angiosperms "receptacle-seed"

Monocots

Unikonta

Opisthokonta

Cnidaria (radial symmetry)

Animalia (Metazoa "beyond animals") Eat other organisms

Arachnids (8 legs)

External Skeleton

Crustaceans

Organs (Eumetazoa "good animals")

Arthropods ("joint-foot")

Insects (6 legs)

| Prokaryotes |
| --- |

Bacteria "stick" (single cell; no nucleus)
Archaea "old" (single cell; no nucleus)
dinoflagellates "terrible whip" [protozoa]
ciliates [protozoa]
brown algae, kelp (Phaeophyceae)
diatoms "cut in two" [algae; plankton]
radiolarians "small sunbeam" [protozoa]
foraminiferans "hole bearers" [plankton]
green algae (Chlorophyta)
mosses (Bryophyta)
liverworts (Marchantiophyta)
ferns (Filicophyta)
cycads (Cycadophyta) [seeds]
conifers (Coniferae) [cones]
rosids
asterids [most flowers]
cacti (Cactaceae)
poppies (Papaveraceae)
laurels (Lauraceae)
magnolias (Magnolia)
lilies (Lilium)
orchids, irises (Asparagales)
palms (Palmae)
grasses (Graminae)
red algae (Rhodophyta) [some seaweeds]
amoeba (Amoebazoa)
slime molds (Mycetozoa)
fungi [yeasts, molds, mushrooms]
comb jellies (Ctenophora "comb bearer")
sponges (Hexactinellid, Calcarea)
corals, anenomes (Anthozoa)
jellyfish (Scyphozoa)
spiders (Araneae)
mites, ticks (Acarina)
scorpions (Scorpiones)
horseshoe crabs (Xiphosura)
Trilobites (extinct) "three lobed"
barnacles (Cirripedia "curl footed")
copepods, krill
crabs, lobster, shrimp (Decapods "ten footed")
millipedes, centipedes (Myriapoda "many feet")
dragonflies (Odonata)
cockroaches (Blattodea)
termites (Isoptera "equal wing")
grasshoppers (Orthoptera "straight wing")
true bugs, cicada, aphid (Hemiptera "half wing")

Plants

Arthropods

This diagram is a *cladogram*, a tree-like picture showing how organisms are related. Each sub-tree in a cladogram is called a *clade*, such as mammals, animals, amphibians. Most branches in a cladogram should split into two sub-trees, but for simplicity this picture has some branches that split into three. Extinct species are represented as dead-end branches. This cladogram is a high-level overview and does not show individual species. Each clade is defined by a distinguishing characteristic that sets it apart from neighboring clades. For example, tetrapods have 4 legs. Sometimes that characteristic disappears in later organisms, for example: snakes are in the tetrapod clade, but no longer have legs. Some well-known groups of organisms are not clades – including reptiles, protists, fish, invertebrates, sponges, and prokaryotes – because they do not include all descendents of the most recent common ancestor.

V3.8
copyright
Neal Olander
tellapallet.com

# Extinctions part of life

Phylogenetic tree showing archosaurs, dinosaurs, birds, etc. through geologic time removed due to copyright restrictions.

# Phylogenetics

**General Problem:**
Infer complete ancestry of
a set of 'objects' based on
knowledge of their 'traits'

'**Objects**' **can be:** Species,
Genes, Cell types, Diseases,
Cancers, Languages, Faiths,
Cars, Architectural Styles

Mammal family tree removed due to copyright restrictions.

'**Traits**' **can be:** Morphological, molecular,
gene expression, TF binding, motifs, words…

**Historical record varies:** Fossils, imprints,
timing of geological events, 'living fossils',
sequencing of extinct species, paintings, stories.

**Today:** Phylogenies using only extant species data
➔ **gene trees** (paralog / ortholog / homolog trees)

# "Central dogma" of Molecular Biology

**DNA** **Epigenomics**

makes

↓

**RNA**

makes

↓

**Protein**

# Chromosomes inside the cell



Figures by MIT OpenCourseWare.

# DNA packaging

- Why packaging
  - DNA is very long
  - Cell is very small
- Compression
  - Chromosome is 50,000 times shorter than extended DNA
- Using the DNA
  - Before a piece of DNA is used for anything, this compact structure must open locally
- Now emerging:
  - Role of accessibility
  - State in chromatin itself
  - Role of 3D interactions

Image removed due to copyright restrictions.
Please see: Figure 8-10 from Alberts, Bruce, and Martin Raff.
Essential Cell Biology. New York, NY: Garland Publishing Inc.,
1997. ISBN: 0815320450.

# Diverse epigenetic modifications



Courtesy of the National Institutes of Health; in the public domain.

Image source: http://nihroadmap.nih.gov/epigenomics/

# Diversity of epigenetic modifications

modifications

Histone tails

DNA wrapped around
histone proteins

- 100+ different histone modifications
  - Histone protein → H3/H4/H2A/H2B
  - AA residue → Lysine4(K4)/K36…
  - Chemical modification → Met/Pho/Ubi
  - Number → Me-Me-Me(me3)
  - Shorthand: H3K4me3, H2BK5ac
- In addition:
  - DNA modifications
  - Methyl-C in CpG / Methyl-Adenosine
  - Nucleosome positioning
  - DNA accessibility
- The constant struggle of gene regulation
  - TF/histone/nucleo/GFs/Chrom compete

# Epigenomics Roadmap across 100+ tissues/cell types



Esophagus
Heart
Aorta
Left ventricle
Right ventricle
Right atrium
Thymus
Lung
Adipose
Breast
Myoepithelial
vHMEC
Progenitor enriched
Luminal epithelial
Duodenum mucosa
Liver
Duodenum Spleen
smooth muscle
Stomach
smooth muscle
Kidney
Pancreas
Small intestine
Psoas muscle
Muscle

Brain
Angular gyrus
Anterior caudate
Cingulate gyrus
Hippocampus middle
Inferior temporal lobe
Substantia Nigra
Dorsolateral
Prefrontal Cortex
Blood
Stem cells (CD34+)
B-Cells (CD19+)
T-Cells (CD3+, CD4+, CD8+)
Granuloytes (CD15+)
PBMCs
NK-Cells (CD56+)
Stomach mucosa
Sigmoid colon
Ovary
Colon
smooth muscle
mucosa
Osteoblasts
Rectum
smooth muscle
mucosa
Germinal matrix

Brain
Thymus
Heart
Aorta
Lung
Right, Left
Cord blood
B-Cells (CD19+)
T-Cells (CD3+)
Liver
Spleen
Placenta

Spinal cord
Stomach
Adrenal
Kidney
Right, Left,
Renal cortex,
Renal pelvis
Small intestine
Large intestine
Skeletal muscle
Back, Trunk, Arm, Leg
Gonad
Ovaries, Testes

Art: Rae Senarighi, Richard Sandstrom

iPS cells
6.9, 18c, 19.11, 20b, 15b
Trophoblast
ES cell lines
H1, H9, I3, WA7, HUES6,
HUES48, HUES64, 4star
Neuronal progenitors
Mesodermal progenitors
Mesenchymal stem cells
Ectoderm
Endoderm

Ganglion Eminence
derived primary
cultured neurospheres
Cortex derived primary
cultured neurospheres

Marrow derived
mesenchymal cells
Chondrocytes

Skin
Skin keratinocyte
Skin fibroblasts
Skin melanocytes

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111
reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

## Diverse epigenomic assays:
1. **Histone modifications**
   - **H3K4me3, H3K4me1**
   - **H3K36me3**
   - **H3K27me3, H3K9me3**
   - **H3K27/9ac, +20 more**
2. **Open chromatin**:
   - **DNase**
3. **DNA methylation**:
   - **WGBS, RRBS, MRE/MeDIP**
4. **Gene expression**
   - **RNA-seq, Exon Arrays**

## Diverse tissues and cells:
1. **Adult tissues and cells** (brain, muscle, heart, digestive, skin, adipose, lung, blood…)
2. **Fetal tissues** (brain, skeletal muscle, heart, digestive, lung, cord blood…)
3. **ES cells, iPS, differentiated cells** (meso/endo/ectoderm, neural, mesench, trophobl)

74

# Deep sampling of 9 reference epigenomes (e.g. IMR90)



Courtesy of Ting Wang. Used with permission. UWash Epigenome Browser, Ting Wang

**Chromatin state+RNA+DNAse+28 histone marks+WGBS+Hi-C** 75

# Diverse chromatin signatures encode epigenomic state

**Enhancers**
- **H3K4me1**
- **H3K27ac**
- **DNase**

**Promoters**
- **H3K4me3**
- **H3K9ac**
- **DNase**

**Transcribed**
- **H3K36me3**
- **H3K79me2**
- **H4K20me1**

**Repressed**
- **H3K9me3**
- **H3K27me3**
- **DNAmethyl**

5'-UTR

3'-UTR

chromatin fiber

DNA

nucleus

nucleosome

- **H3K4me3**
- **H3K4me1**
- **H3K27ac**
- **H3K36me3**
- **H4K20me1**
- **H3K79me3**
- **H3K27me3**
- **H3K9me3**
- **H3K9ac**
- **H3K18ac**

- 100s of known modifications, many new still emerging
- Systematic mapping using ChIP-, Bisulfite-, DNase-Seq

# Chromatin state annotations across 127 epigenomes



| | Emissions |
|---|---|
| 1 Active TSS | |
| 2 Flanking Active TSS | |
| 3 Transcr. at gene 5' and 3' | |
| 4 Strong transcription | |
| 5 Weak transcription | |
| 6 Genic enhancers | |
| 7 Enhancers | |
| 8 ZNF genes & repeats | |
| 9 Heterochromatin | |
| 10 Bivalent/Poised TSS | |
| 11 Flanking Bivalent TSS/Enh | |
| 12 Bivalent Enhancer | |
| 13 Repressed Polycomb | |
| 14 Weak Repressed Polycomb | |
| 15 Quiescent/Low | |

H3K4me3  H3K4me1  H3K36me3  H3K9me3  H3K27me3

Courtesy of Anshul Kundaje. Used with permission.

Reveal epigenomic variability: enh/prom/tx/repr/het

Anshul Kundaje

# "Central dogma" of Molecular Biology

# Genes control the making of cell parts

- The gene is a fundamental unit of inheritance
  - Each DNA molecule ⇔ 10,000+ genes
  - 1 gene ⇔ 1 functional element (one "part" of cell machinery)
  - Every time a "part" is made, the corresponding gene is:
    - Copied into mRNA, transported, used as blueprint to make protein
- RNA is a temporary copy
  - The medium for transporting genetic information from the DNA information repository to the protein-making machinery is an RNA molecule
  - The more parts are needed, the more copies are made
  - Each mRNA only lasts a limited time before degradation

# mRNA:  The messenger



Image by MIT OpenCourseWare.

- Information changes medium
  - single strand vs. double strand
  - ribose vs. deoxyribose sugar

| A | T | T | A | C | G | G | T | A | C | C | G | T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| U | A | A | U | G | C | C | A | U | G | G | C | A |

  - Compatible base-pairing in hybrid



uracil (RNA)          thymine (DNA)

# From DNA to RNA: Transcription

Image removed due to copyright restrictions.Please see: Figure 7-9 from Alberts, Bruce, and Martin Raff. Essential Cell Biology. New York, NY: Garland Publishing Inc., 1997. ISBN: 0815320450.

# From pre-mRNA to mRNA:  Splicing

- In Eukaryotes, not every part of a gene is coding
  - Functional exons interrupted by non-translated introns
  - During pre-mRNA maturation, introns are spliced out
  - In humans, primary transcript can be $10^6$ bp long

Image removed due to copyright restrictions.
Please see: Figure 7-16 from Alberts, Bruce, and Martin Raff. Essential Cell Biology.
New York, NY: Garland Publishing Inc., 1997. ISBN: 0815320450.

  - Alternative splicing can yield different exon subsets for the same gene, and hence different protein products

# RNA can be functional

- Single Strand allows complex structure
    - Self-complementary regions form helical stems
    - Three-dimensional structure allows functionality of RNA
- Four types of RNA
    - mRNA: messenger of genetic information
    - tRNA: codon-to-amino acid specificity
    - rRNA: core of the ribosome
    - snRNA:  splicing reactions
- To be continued…
    - We'll learn more in a dedicated lecture on RNA world
    - Once upon a time, before DNA and protein, RNA did all

# RNA structure:  2$^{nd}$ary and 3$^{rd}$ary



Courtesy of SStructView

84

# Splicing machinery made of RNA

Image removed due to copyright restrictions.
Please see: Figure 7-16 from Alberts, Bruce, and Martin Raff. Essential Cell Biology.
New York, NY: Garland Publishing Inc., 1997. ISBN: 0815320450.

# "Central dogma" of Molecular Biology

**DNA**

↓ makes

**RNA**

↓ makes

**Protein**

# Proteins carry out the cell's chemistry



DNA

Replication

Transcription

RNA

Translation

Protein

Image by MIT OpenCourseWare.

- **More complex polymer**
  - Nucleic Acids have 4 building blocks
  - Proteins have 20. Greater versatility
  - Each amino acid has specific properties

**Sequence → Structure → Function**
  - The amino acid sequence determines the three-dimensional fold of protein
  - The protein's function largely depends on the features of the 3D structure

- **Proteins play diverse roles**
  - Catalysis, binding, cell structure, signaling, transport, metabolism

# Protein structure



Sugar phosphate backbone

DNA

A

2

3

1

Base pair

## Helix-turn-helix

Common motif for DNA-binding proteins that often play a regulatory role as mRNA level transcription factors

## Beta-barrel

Some antiparallel b-sheet domains are better described as b-barrels rather than b-sandwiches, for example streptavadin and porin. Note that some structures are

intermediate between the extreme barrel and sandwich arrangements.

## Alpha-beta horseshoe

this placental ribonuclease inhibitor is a cytosolic protein that binds extremely strongly to any ribonuclease that may leak into the cytosol. 17-stranded parallel b sheet curved into an open horseshoe shape, with 16 a-helices packed against the outer surface. It doesn't form a barrel although it looks as though it should. The strands are only very slightly slanted, being nearly parallel to the central `axis'.

# Protein building blocks

- Amino Acids

# From RNA to protein: Translation



Image by MIT OpenCourseWare.

# The Genetic Code



→ Use evolutionary and compositional properties
to computationally discover protein-coding genes

# Summary: The Central Dogma

## DNA makes RNA makes Protein

Inheritance

Messages

Reactions

DNA

Replication

Transcription

RNA

Translation

Protein

Image by MIT OpenCourseWare.

# Cellular dynamics and regulation
## *How cells move through this Central Dogma*



**DNA**

makes

**Gene regulation**

**RNA**

makes

**Protein**

# Animal/Human gene regulation:
# One genome ⇔ Many cell types

ACCAGTTACGACGGTCA
GGGTACTGATACCCCAA
ACCGTTGACCGCATTTA
CAGACGGGGTTTGGGTT
TTGCCCCACACAGGTAC
GTTAGCTACTGGTTTAG
CAATTTACCGTTACAAC
GTTTACAGGGTTACGGT
TGGGATTTGAAAAAAG
TTTGAGTTGGTTTTTTC
ACGGTAGAACGTACCGT
TACCAGTA



Image in the public domain.

Images of a heart, red blood cell, and a brain removed due to copyright restrictions.

# Eukaryotic Gene Regulation

Cartoon depicting eukaryotic gene regulation removed due to copyright restrictions.

# Diverse roles for regulatory non-coding RNAs

- **Small RNA pathways (18-21 nt)**
  - microRNAs:
    - Repress genes by targeting their 3′ UTRs by complementarity
    - Double-stranded RNA is then recognized and degraded
    - Recently found to also target promoter regions in rare cases
  - piwiRNAs
    - Target and repress transposable elements in germline
  - snoRNAs
  - 21U-RNAs
- **Long non-coding RNAs (1000s nt, many exons)**
  - Scaffolds for protein/TF binding
  - Scaffolds for 3D structure of RNA

# Regulation of Gene Expression



**Transcription Factor**

**Polymerase**

**Promoter**

mRNA

**Transcription Factor Binding Site**

Examples:

- Upstream of genes are *promoter* regions
- Contain promoter sequences or *motifs*
- *Transcription factors* (TFs) bind to motifs
- TFs recruit *RNA polymerase*
- Gene transcription

# Predicted motif drivers of enhancer modules



- **Activator and repressor motifs consistent with tissues**

**Pouya Kheradpour**

Source: Zeitlinger, Julia et al. "Whole-genome ChIP–chip analysis of Dorsal, Twist, and Snail suggests integration of diverse patterning processes in the Drosophila embryo." Genes & Development 21, no. 4 (2007): 385-390.

- Feed-forward loops in developmental patterning
- Cooperation of master reg. & downstream reg.

# Systematic motif dissection in 2000 enhancers: 5 activators and 2 repressors in 2 cell lines

Figure 1: selection of activator and repressor motifs removed due to copyright restrictions.
Source: Kheradpour, Pouya et al. "Systematic dissection of regulatory motifs in 2000 predicted human enhancers using a massively parallel reporter assay." Genome Research 23, no. 5 (2013): 800-811.

**54000+ measurements (x2 cells, 2x repl)**

**Kheradpour *et al* <u>Genome Research</u> 2013**

# Emerging properties of regulatory networks

Figures removed due to copyright restrictions.

- Hierarchical levels of regulatory control
    - Small number of backward-pointing edges
- Specific / distinct feedback by microRNAs at each level
    - Two classes of TFs: miRNA regulators and miR-regulated

# From Systems Biology to Synthetic Biology

**Synthetic Regulatory Networks**

| Design & model network | Encode into DNA plasmid | Transfer to cell | Test network dynamics |
|---|---|---|---|

**Synthetic Metabolic Pathways**

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Benner, Steven A. and A. Michael Sismour. "Synthetic biology." Nature Reviews Genetics 6, no. 7 (2005): 533-543.

**Jim Collins**

- Components with known properties
- Assemble based on engineering goals / principles
- Implement within engineered cells and organisms
- Study behavior & adjust as needed

**Jay Keasling**

# Over-express a single microRNA leads to new wing



A — →wing, haltere — Bx>mir-iab-4 anti-sense w1118
B — B' Sensory bristles
C — haltere — WT — w1118
D — →wing — sense — Bx>mir-iab-4 sense
E — →wing w/bristles — Antisense — Bx>mir-iab-4 anti-sense

Note: C,D,E same magnification

- Discovery of sense/anti-sense miRNAs
- Regulatory switch selects between two developmental programs
- By over-expressing one strand (miRNAas) the balance is tilted
- Wing program launched vs. haltere

**Stark _et al_, Genes&Development 2007**

# Brief intro to Human Genetics

# The role of genetic alterations

# Brief intro to human genetics

- **Human genome**: 3.2B letters, 2 copies, 23 chromosomes, 20k genes, ~3M common SNPs, ~500k haplotype blocks
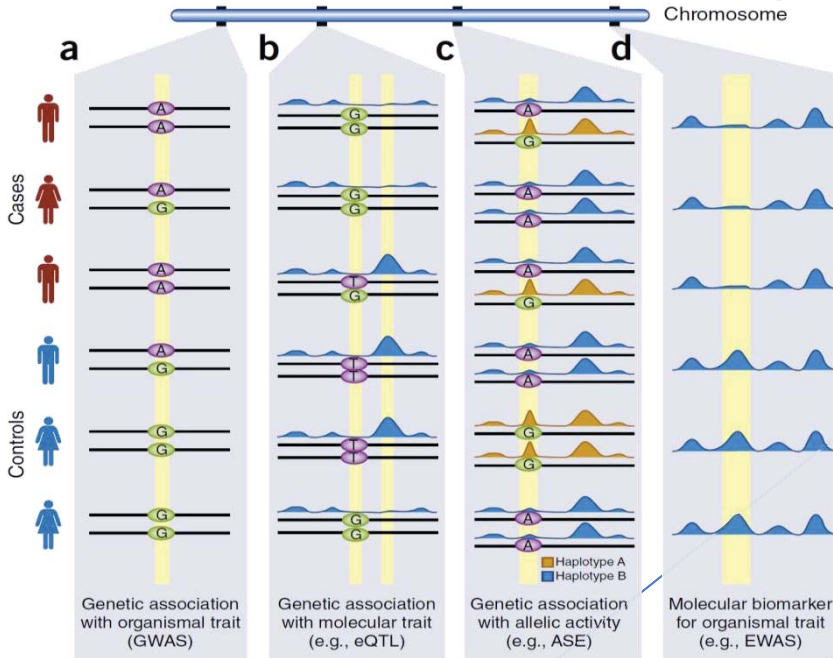


Figure in the public domain. Created by Darryl Leja and Teri Manolio, NHGRI; Tony Burdett, Dani Welter, and Helen Parkinson, EBI.

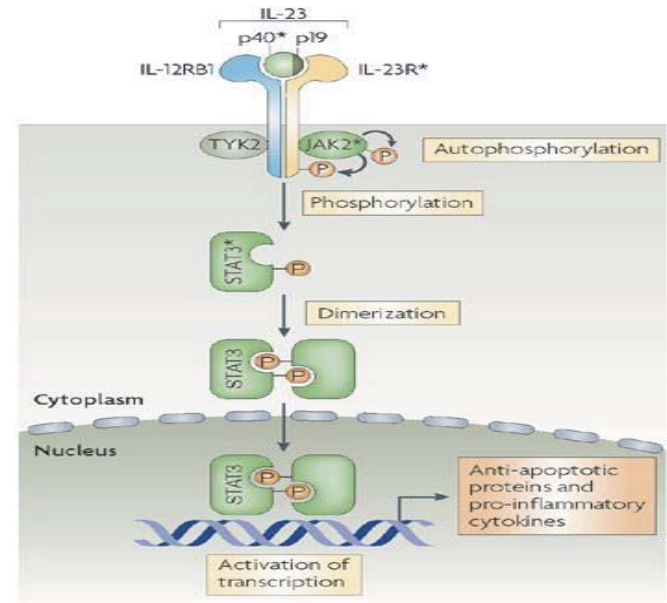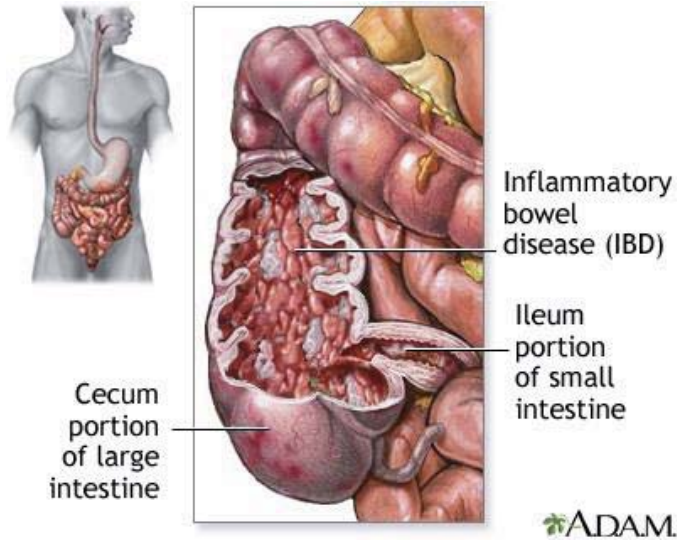# The power and challenge of disease-association studies

Slide credit: Luke Ward, Mark Daly

- **Large associated blocks with many variants: Fine-mapping challenge**
- **No information on cell type/mechanism, most variants non-coding**
- ➔ **Epigenomic annotations help find relevant cell types / nucleotides** 107
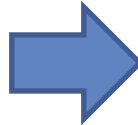
# The power of GWAS: reveal new disease genes



Inflammatory bowel disease (IBD)

Ileum portion of small intestine

Cecum portion of large intestine

*ADAM.*

Nature Reviews | Immunology

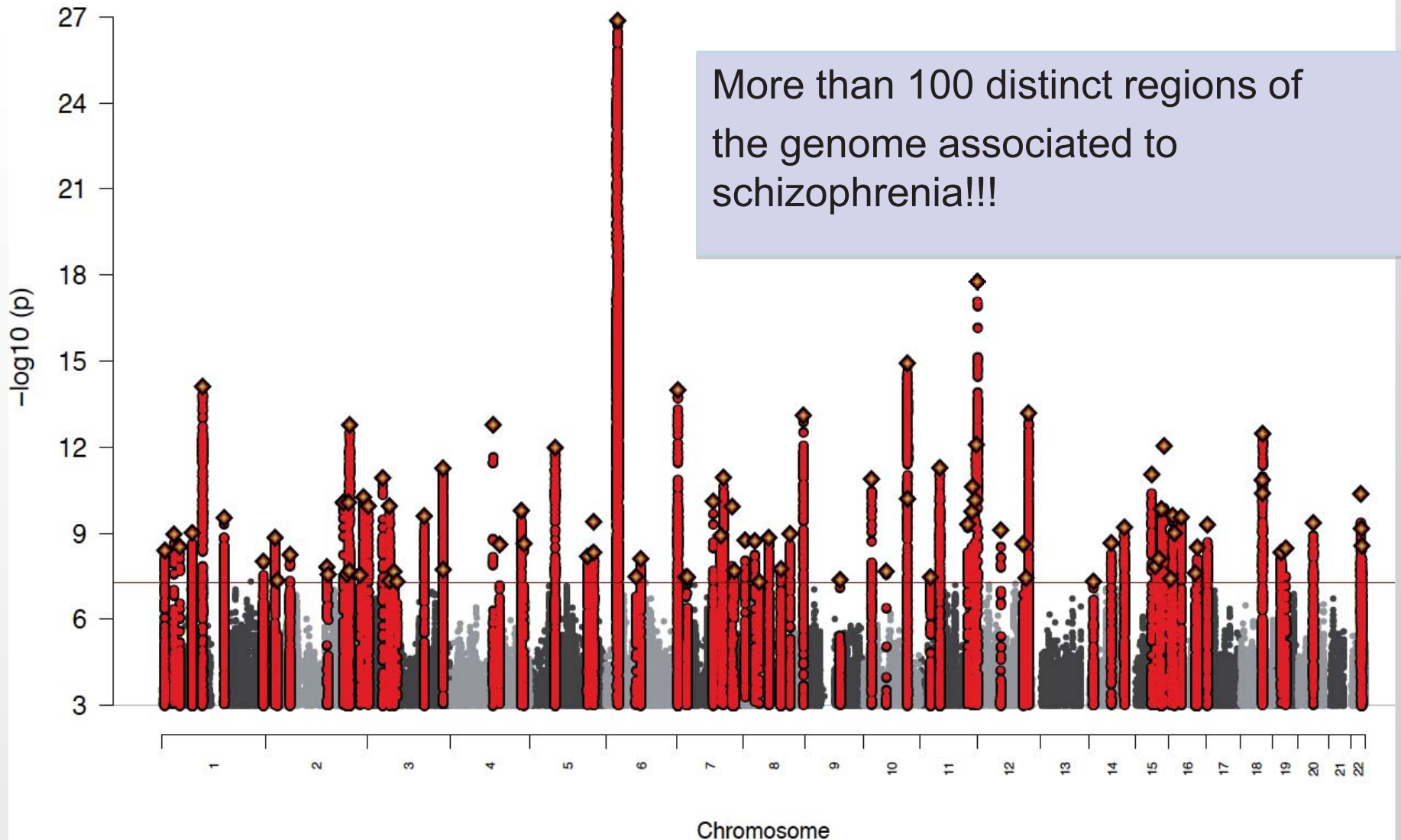| rs11209026 | A | G |
|------------|-----|-----|
| Cases | 22 | 976 |
| Controls | 68 | 932 |

**Chi-sq = 24.5, p=7.3 x $10^{-7}$**

IL23R cytokine receptor on a subset of effector T-cells

# Genomewide association in schizophrenia with 40,000 cases



More than 100 distinct regions of the genome associated to schizophrenia!!!

Source: Ripke, Stephan et al. "Biological insights from 108 schizophrenia-associated genetic loci." Nature 511, no. 7510 (2014): 421.

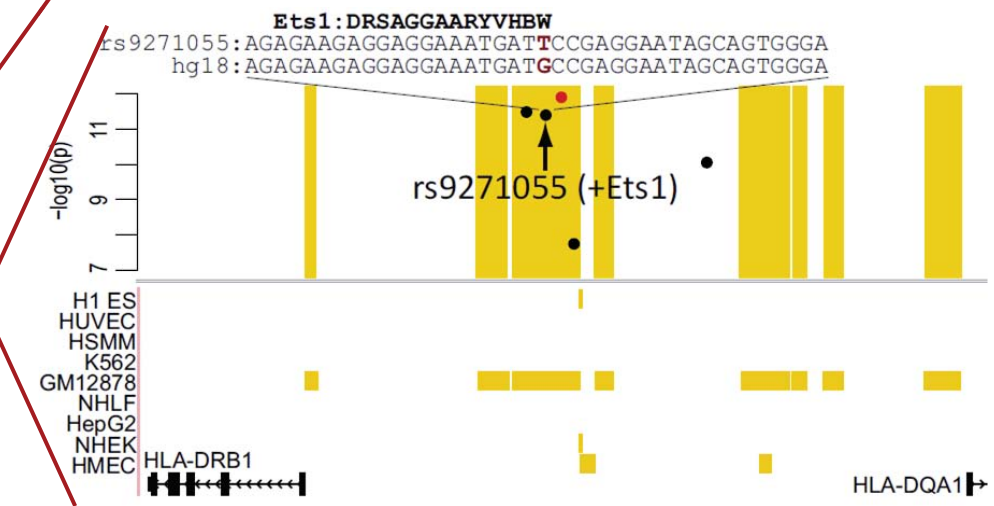# Interpreting non-coding variants

| Phenotype | Top Cell Type | Total #SNPs from Study | #SNPs in enh. States 4 and 5 | p-value | FDR | H1 ES | K562 | GM12878 | HepG2 | HUVEC | HSMM | NHLF | NHEK | HMEC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Erythrocyte phenotypes (Ref. 38) | K562 | 35 | 9 | $<10^{-7}$ | 0.02 | 9 | 17 | 4 | 0 | 0 | 1 | 2 | 1 | 1 |
| Blood lipids (Ref. 39) | HepG2 | 101 | 13 | $<10^{-7}$ | 0.02 | 3 | 5 | 0 | 11 | 2 | 3 | 3 | 4 | 3 |
| Rheumatoid arthritis (Ref. 40) | GM12878 | 29 | 7 | $2.0 \times 10^{-7}$ | 0.03 | 0 | 0 | 15 | 0 | 2 | 0 | 0 | 2 | 3 |
| Primary billary cirrhosis (Ref. 41) | GM12878 | 6 | 4 | $6.0 \times 10^{-7}$ | 0.03 | 0 | 11 | 41 | 0 | 0 | 0 | 0 | 8 | 8 |
| Systemic lupus erythromatosus (Ref. 42) | GM12878 | 18 | 6 | $9.0 \times 10^{-7}$ | 0.03 | 0 | 4 | 21 | 0 | 5 | 8 | 0 | 3 | 5 |
| Lipoprotein cholesterol/triglycerides (Ref. 43) | HepG2 | 18 | 5 | $1.2 \times 10^{-6}$ | 0.03 | 17 | 8 | 0 | 24 | 3 | 6 | 4 | 3 | 3 |
| Hematological traits (Ref. 44) | K562 | 39 | 7 | $1.7 \times 10^{-6}$ | 0.03 | 0 | 12 | 10 | 2 | 1 | 0 | 0 | 1 | 0 |
| Hematological parameters (Ref. 45) | K562 | 28 | 6 | $2.2 \times 10^{-6}$ | 0.03 | 0 | 15 | 7 | 0 | 5 | 7 | 7 | 3 | 2 |
| Colorectal cancer (Ref. 46) | HepG2 | 4 | 3 | $3.8 \times 10^{-6}$ | 0.03 | 0 | 0 | 0 | 66 | 0 | 12 | 0 | 12 | 12 |
| Blood pressure (Ref. 47) | K562 | 9 | 4 | $5.0 \times 10^{-6}$ | 0.04 | 0 | 30 | 14 | 0 | 10 | 6 | 7 | 5 | 11 |



- **Disease-associated SNPs enriched for enhancers in relevant cell types**
- **E.g. lupus SNP in GM enhancer disrupts Ets1 predicted activator**

# Mechanistic predictions for top disease-associated SNPs

**Lupus erythromatosus in GM lymphoblastoid**    **Erythrocyte phenotypes in K562 leukemia cells**

Figures removed due to copyright restrictions.

**Disrupt activator Ets-1 motif**
➔ **Loss of GM-specific activation**
➔ **Loss of enhancer function**
➔ **Loss of HLA-DRB1 expression**

**Creation of repressor Gfi1 motif**
➔ **Gain K562-specific repression**
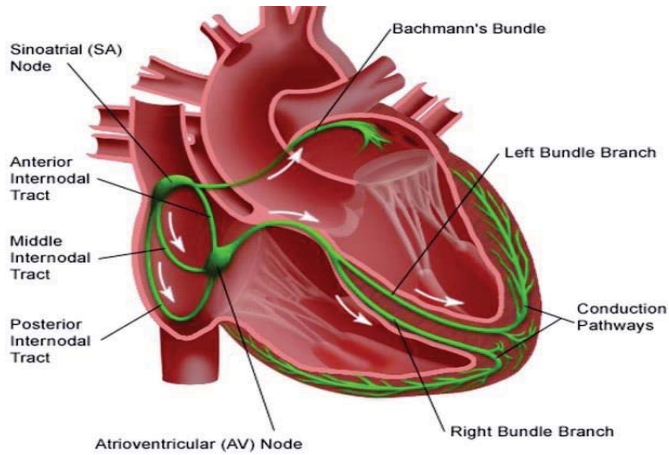➔ **Loss of enhancer function**
➔ **Loss of CCDC162 expression**

# Chromatin state annotations across 127 epigenomes

Figures removed due to copyright restrictions.

# Reveal epigenomic variability: enh/prom/tx/repr/het

Anshul Kundaje

# Characterizing sub-threshold variants in heart arrhythmia

*Focus on sub-threshold variants*
*(e.g. rs1743292 P=10$^{-4.2}$)*

Courtesy of Macmillan Publishers Limited. Used with permission.
Source: Arking, Dan E. et al. "Genetic association study of QT interval highlights role for calcium signaling pathways in myocardial repolarization." Nature Genetics 46, no. 8 (2014): 826-836.

*Trait: QRS/QT interval*
(1) Large cohorts, (2) many known hits
(3) well-characterized tissue drivers

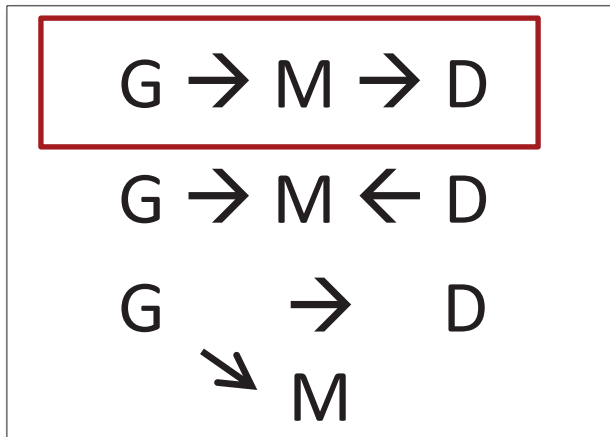| Trait | Most enriched tissue/cell type (Abbrev) | −logP |
|---|---|---|
| Height | ESC | 4.7 |
| Height | ESC | 4.0 |
| Crohn's disease | Tper | 7.7 |
| Chronic lymphocytic leukaemia | Tcor | 4.9 |
| Type 1 diabetes autoantibodies | Treg | 4.6 |
| Type 1 diabetes | Treg | 4.1 |
| Platelet counts | Th.nai | 4.6 |
| Chronic lymphocytic leukaemia | Th.stm | 5.7 |
| Self-reported allergy | Th.stm | 4.9 |
| Graves' disease | Th.stm | 4.3 |
| Celiac disease | Th17st | 6.9 |
| Rheumatoid arthritis | Th17st | 4.2 |
| Multiple sclerosis | Th.mm | 11.6 |
| Celiac disease + rheum. arthritis | Th.mm | 5.6 |
| Type 1 diabetes | Th.mm | 5.5 |
| Systemic lupus erythematosus | Th.mm | 4.8 |
| Systemic lupus erythematosus | Bcor | 5.4 |
| Primary biliary cirrhosis | Bcor | 3.9 |
| Red blood cell traits | HSCmb | 5.9 |
| Platelet counts | HSCmb | 8.0 |
| Mean platelet volume | HSCmb | 5.0 |
| Mean platelet volume | HSCmb | 3.9 |
| Rheumatoid arthritis | Bper | 8.5 |
| Multiple sclerosis | Bper | 4.7 |
| Rheumatoid arthritis | NKper | 5.0 |
| Mean platelet volume | Fat | 4.2 |
| HDL cholesterol | Fat | 4.9 |
| Height | Fblast | 4.8 |
| Multiple myeloma | Thym | 4.2 |
| Adiponectin levels | Brain | 4.3 |
| Attention deficit hyperact. disord. | Brain | 4.5 |
| PR interval | Heart | 4.7 |
| Blood pressure | Heart | 4.5 |
| Aortic root size | Vascl | 4.1 |
| Pulmonary function | SmMu | 4.2 |
| Liver enzyme levels (g-glut tx) | Gl.Int | 4.9 |
| Urate levels | Gl.Int | 4.5 |
| Adv. resp. to chemth. (neutr/leuc) | Gl.Muc | 4.0 |
| Breast cancer | Stomc | 4.5 |
| Type 2 diabetes | Stomc | 4.3 |
| Insulin-like growth factors | Placnt | 4.2 |
| Fasting glucose-related traits | P.islets | 4.1 |
| LDL cholesterol | Liver | 10.1 |
| Cholesterol, total | Liver | 9.0 |
| Cholesterol, total | Liver | 7.1 |
| LDL cholesterol | Liver | 6.8 |
| Lipid metabolism phenotypes | Liver | 5.8 |
| HDL cholesterol | Liver | 5.7 |
| Cholesterol, total | Liver | 4.8 |
| HDL cholesterol | Liver | 3.9 |
| Metabolite levels | Liver | 3.9 |
| Platelet counts | T.Leuk | 4.5 |
| Primary biliary cirrhosis | Lymph | 6.7 |
| Mean corpuscular volume | Leuk | 4.7 |
| Inflammatory bowel disease | Mncyt | 14.6 |
| Ulcerative colitis | Mncyt | 6.3 |
| Alzheimer's disease (late onset) | Mncyt | 4.9 |
| Pre-eclampsia | Bone | 4.5 |

Courtesy of Macmillan Publishers Limited. Used with permission.

Source: Roadmap Epigenomics Consortium et al. "Integrative analysis of 111 reference human epigenomes." Nature 518, no. 7539 (2015): 317-330.

# Linking traits to their relevant cell/tissue types

# Methylation differences a causal component of AD



AD-assoc. haplotypes vs. AD-assoc. probes

15k probes

*Methylation probes altered in AD are enriched in AD-associated SNPs*



G → M → D

G → M ← D

G → D
↘
M

*Set-wise causality testing*



AD classification accuracy (AUC)

APOE4 Dosage — 0.59
SNPs — All SNPs 0.59, meQTLs 0.53 / 0.60
Covariates — SVs 0.55, +KC 0.61 / 0.69
All Meth. — 0.71 / 0.73
meQTL Targ. — 0.776, 0.75 / 0.73, meQTL effect removed 0.72
EnhA1 Meth. — 0.74 / 0.77
TssA Meth. — 0.67 / 0.69
Het. Meth. — 0.60 / 0.68

Increase using APOE4 dosage

+ Covariates

*AD predictive power reduced after removing meQTL effect*

116

# Uncovering the molecular basis of top obesity gene

**Lean**

**Obese**

**ARID5B KD
(obesity)**

**ARID5B OE
(anti-obesity)**

**IRX3, IRX5 knock-down ★
(anti-obesity phenotypes)**

**IRX3, IRX5 overexpression
(pro-obesity phenotypes)**

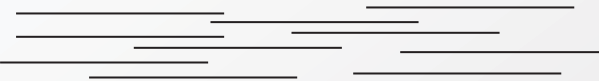**★ C-to-T motif rescue
(anti-obesity phenotypes)**

**T-to-C motif disruption
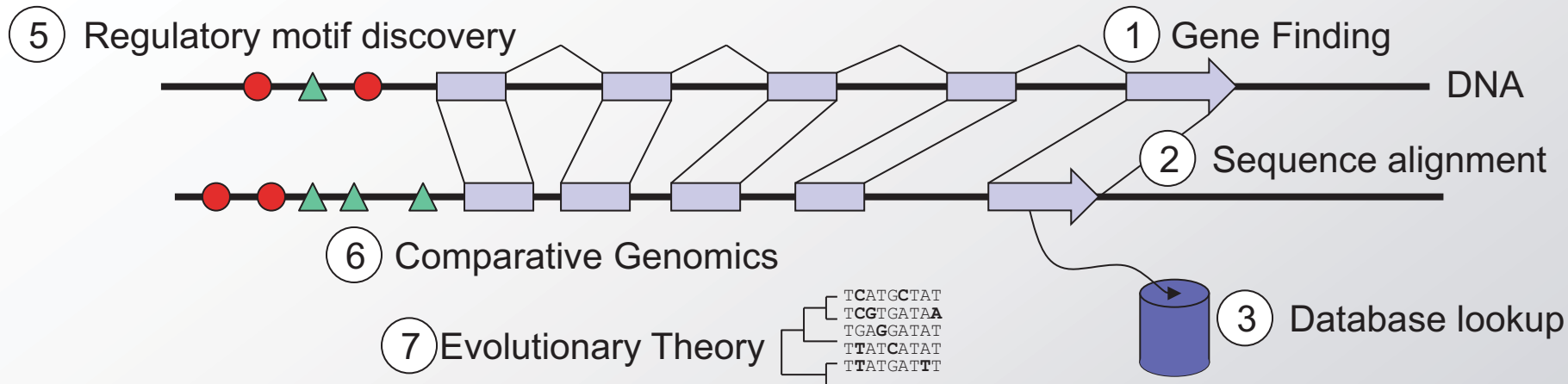(pro-obesity phenotypes)**

# Model: beige ⇔ white adipocyte development

*Shift therapeutic focus from brain to adipocytes*

# Challenges in Computational Biology

④ Genome Assembly

⑤ Regulatory motif discovery

① Gene Finding

DNA

② Sequence alignment

⑥ Comparative Genomics

⑦ Evolutionary Theory

```
TCATGCTAT
TCGTGATAA
TGAGGATAT
TTATCATAT
TTATGATTT
```

③ Database lookup

⑧ Gene expression analysis

RNA transcript

⑨ Cluster discovery

⑩ Gibbs sampling

⑪ Protein network analysis

⑫ Metabolic modelling

⑬ Emerging network properties

6.047 / 6.878 / HST.507 Computational Biology
Fall 2015