All Life depends on 3 critical molecules DNAs (Deoxyribonucleicsield)

- Hold information on how sel works preview page 6
- RNAs (Ribonucleic acid)
 - Act to transfer short pieces of information to different parts of cell
 - Provide templates to synthesize into protein
- Proteins
 - Form enzymes that send signals to other cells and regulate gene activity
 - Form body's major components



This process is called *splicing*



Problems

- Sanger sequencing error rate 6 ease varies from 1% to 3%¹ Repeats in DNA
- ample, ~30000 se longs Alu sequence repeated is over million times in
 - Repeats occur in different scales
- What happens if repeat length is longer than read length? ۲
- Shortest superstring problem ۲
 - Find the shortest string that "explains" the reads —
 - Given a set of strings (reads), find a shortest string that contains all of them

Repeats in DNA and genome assembly CO.UK CO.UK Co.UK Co.UK Co.UK Co.UK Co.UK 00 rpt1B III Ш

Figure 2. Repeat sequence. The top represents the correct layout of three DNA sequences. The bottom shows a repeat collapsed in a misassembly.

BAC-by-BAC sequencing

- Each BAC (Bacterial Artificial Artificial
- Covering the human genome requires ~30000 BACs
- BACs congun-sequenced separately
- Number of repeats in each BAC is significantly smaller than in the whole genome...
 - ...needs much more manual work compared to whole-genome shotgun sequencing

Sequencing of the Human Genome

- The (draft) human genome was ale . Co.uk published in 2001 NO 52 from 45 of 52 Human •
- - Human Genome Project (public consortium)
 - Celera (private company)
- HGP: BAC-by-BAC approach ullet
- Celera: whole-genome shotgun ulletsequencing





HGP: Nature 15 February 2001 Vol 409 Number 6822

> Celera: Science 16 February 2001 Vol 291, Issue 5507





• Indicates positions of interrogation Ligation cycle **2** 2 3 4 5 6 **2** Mardis ER. 2008. Annu. Rev. Genomics Hum. Genet. 9:387–402

a primer to the shared adapter sequences on each amplified fragment, and then DNA ligase is provided along with specific fluorescent-labeled 8mers, whose 4th and 5th bases are encoded by the attached fluorescent group. Each ligation step is followed by fluorescence detection, after which a regeneration step removes bases from the ligated 8mer (including the fluorescent group) and concomitantly prepares the extended primer for another round of ligation. (*b*) Principles of twobase encoding. Because each fluorescent group on a ligated 8mer identifies a two-base combination, the resulting sequence reads can be screened for basecalling errors versus true polymorphisms versus single base deletions by aligning

the individual reads to a known high-quality reference sequence.