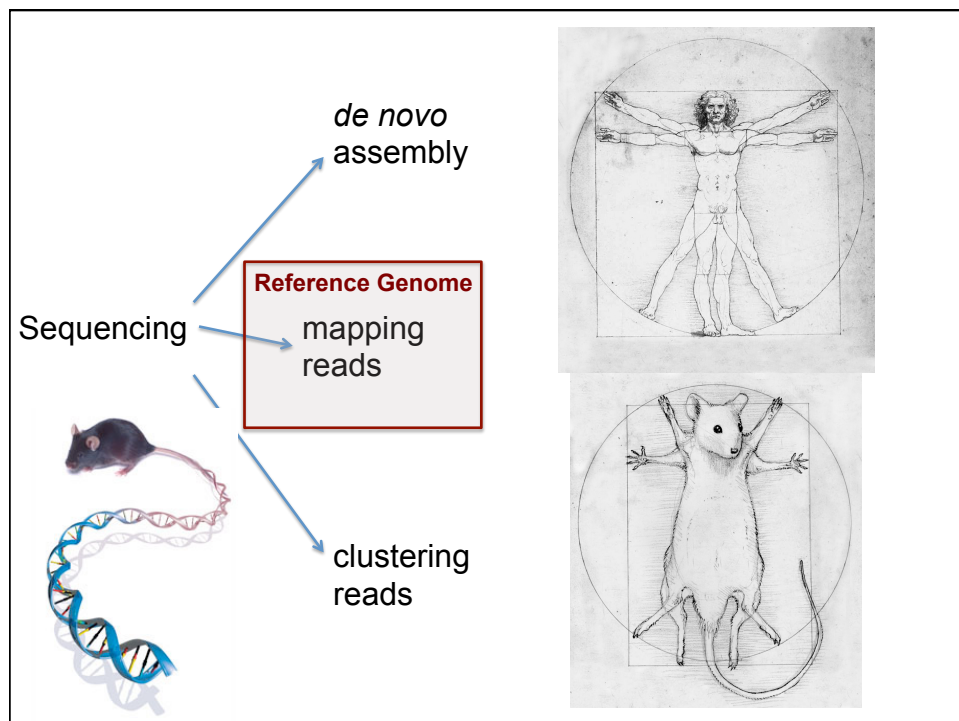
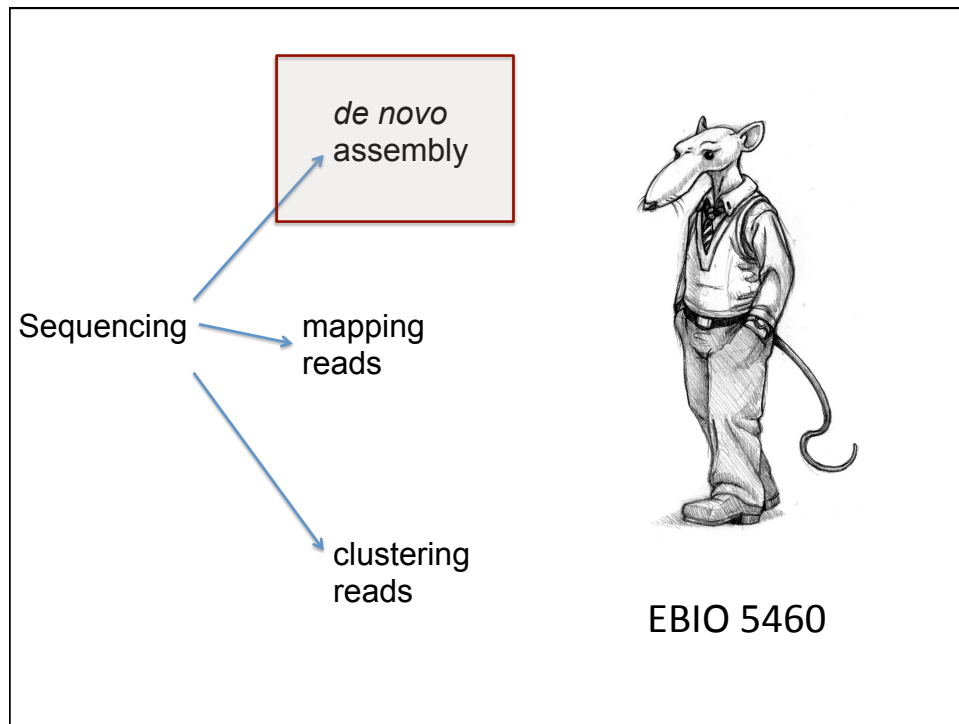


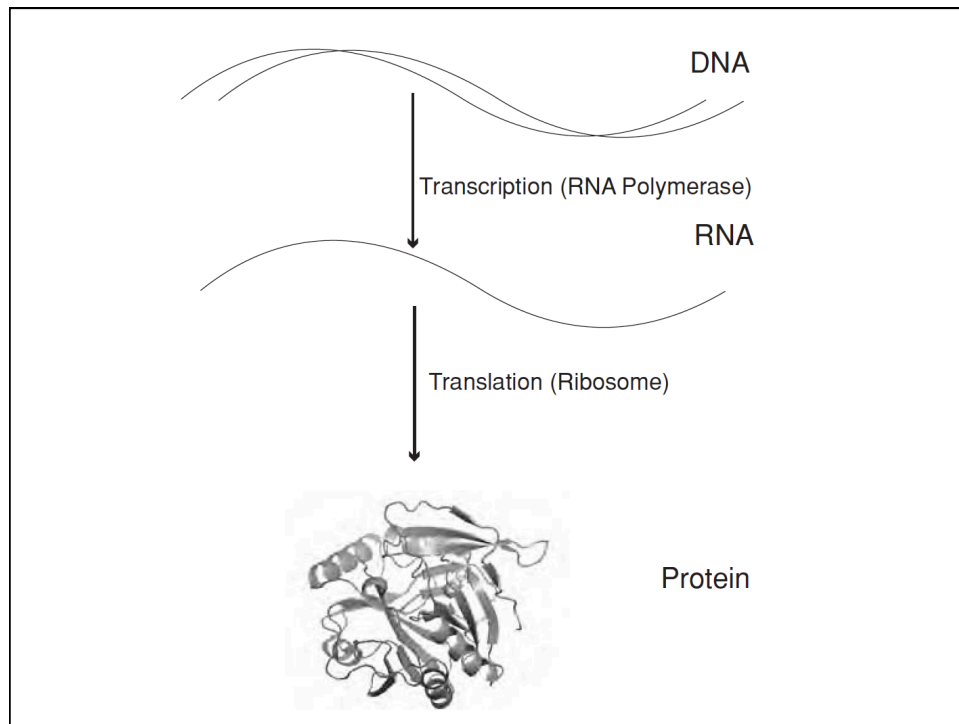
"The availability of genome sequence is just the beginning. Scientists now want to understand the genes and the role they play in the prevention, diagnosis and treatment of disease."

Dr Randy Scott, President of Incyte

Final words on assembly

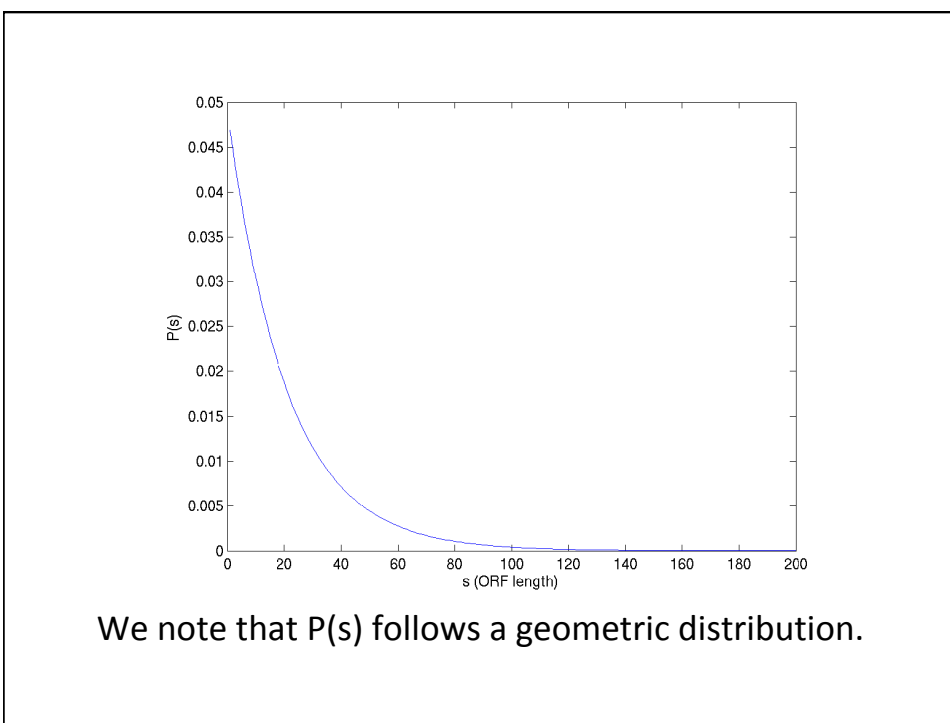
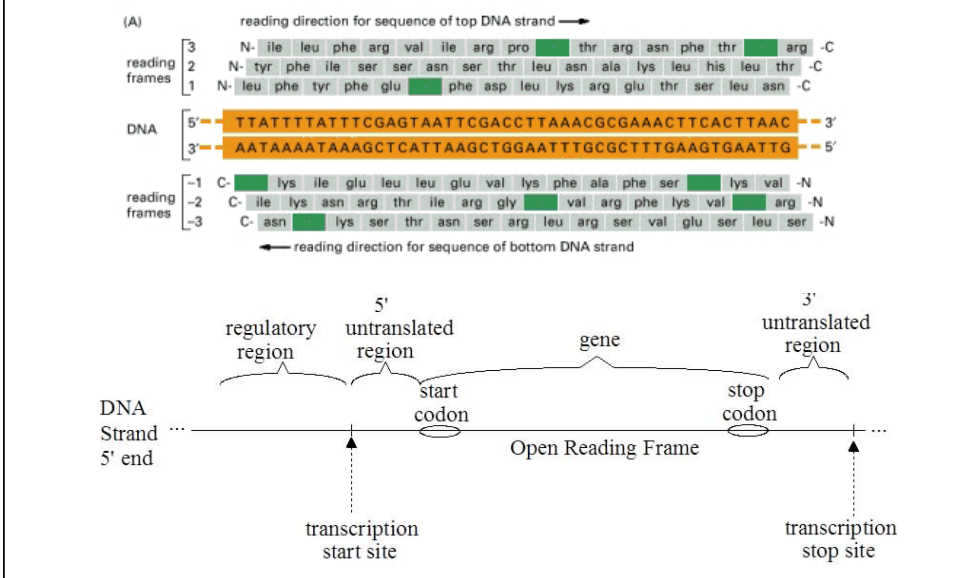
- Assembly is resource intensive. As genomes get bigger, need extensive processing power (compute clusters).
- Repeats increase the time and space of these algorithms exponentially.
- Errors in the fragments from the sequencing instruments confound assembly.

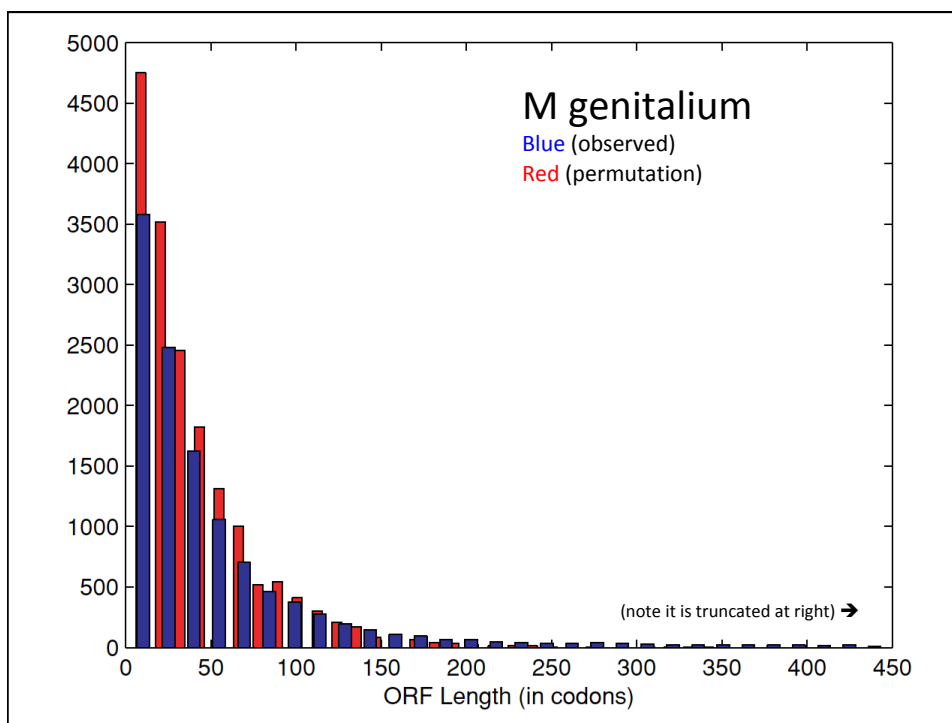




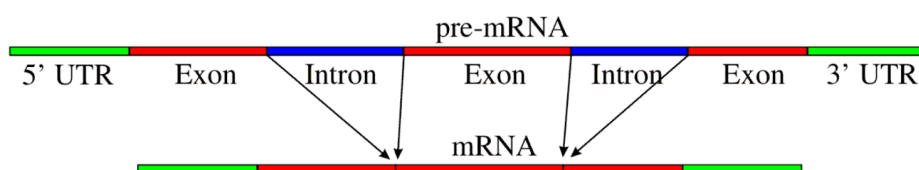
		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA Stop UAG Stop	UGU } Cys UGC } UGA Stop UGG Trp	U C A G
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G
	A	AUU } AUC } Ile AUA } AUG Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }	U C A G

Any region of the DNA sequence can, in principle, code for six different amino acid sequences, because any one of three different reading frames can be used to interpret each of the two strands.





Eukaryotic gene finding (harder)



- Statistical models of codon usage
- Markov models of gene structure
- Comparative genomics approaches