

# The challenges of big data in Genomics

*Antonio Barbadilla*

*Group Genomics, Bioinformatics & Evolution*

*Institut Biotecnologia i Biomedicina*

*Departament de Genètica i Microbiologia*

*UAB*



## Outline

Genome science: the HGP, *a new starting point*

---

The essence of Genomics

---

Genome sequencing

---

Steps of genome analysis

---

The technological explosion: Genome sequences as commodity

---

The triumphal march of genomics

---

Genome science challenges

---

The foreseeable future of genome science

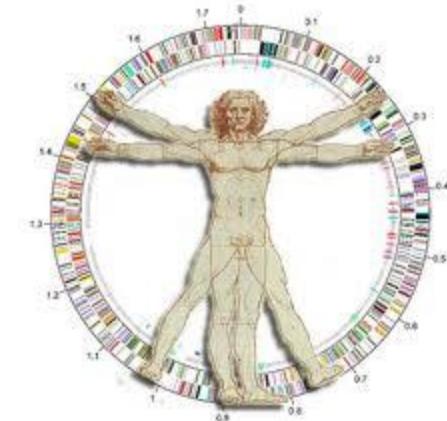
---

Readings

---



# Genome science: the HGP, *a new starting point*

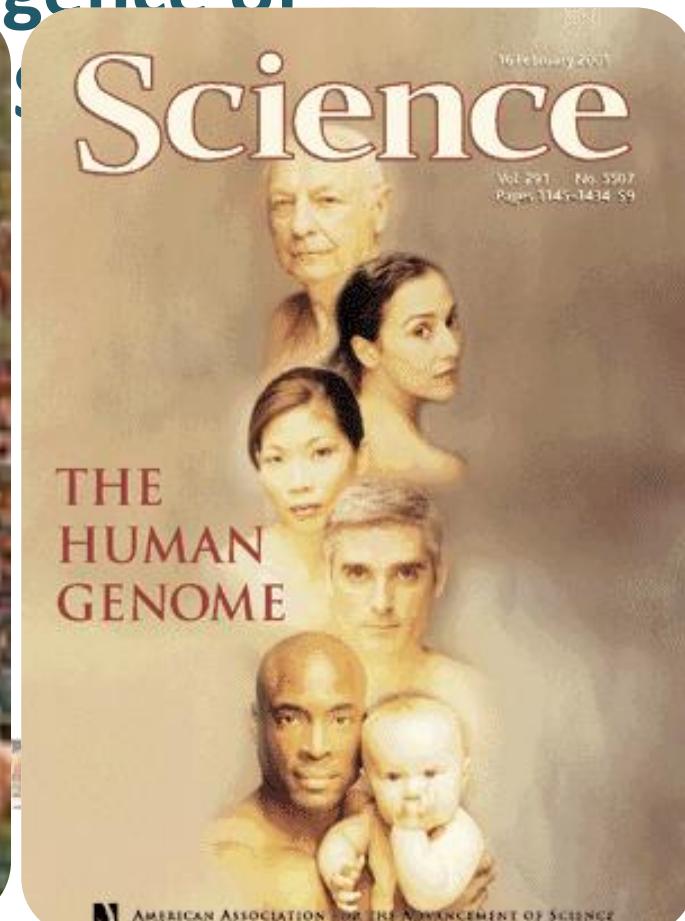
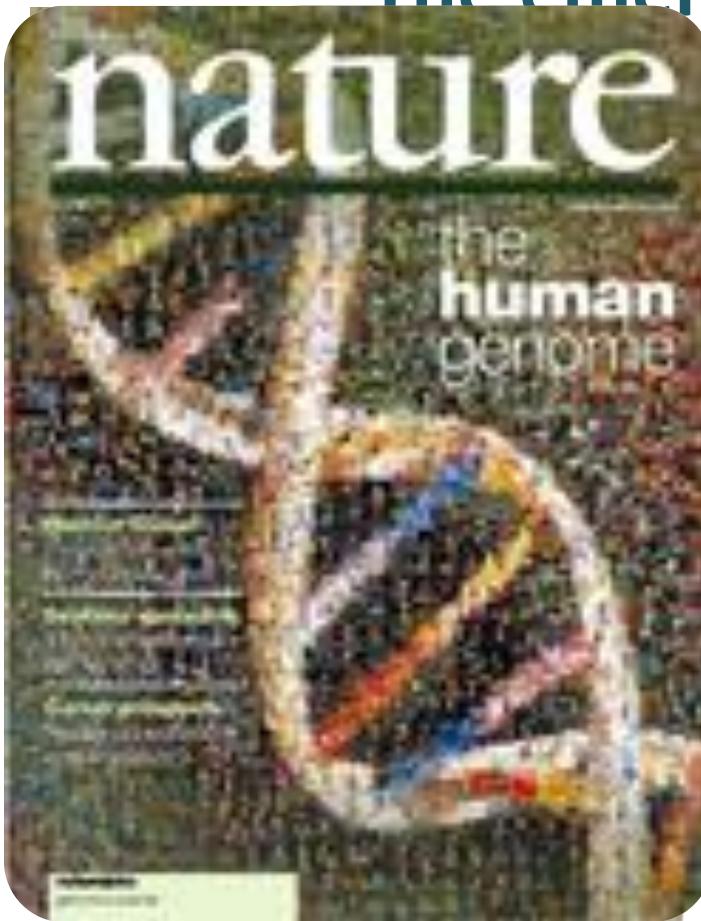




# Genomics Landmarks 2001 and 2004

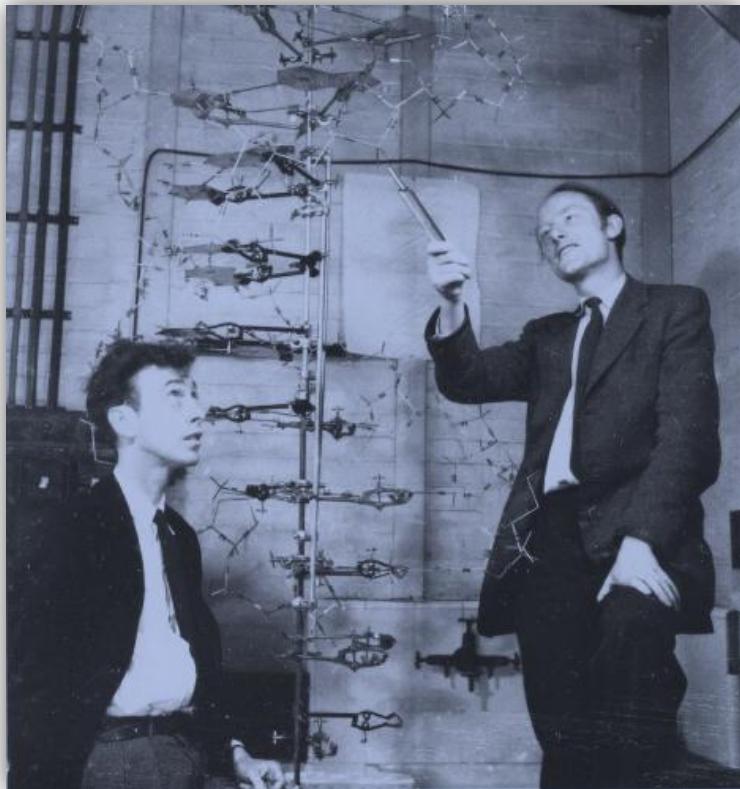
The publication of the draft and complete human sequence

The emergence of

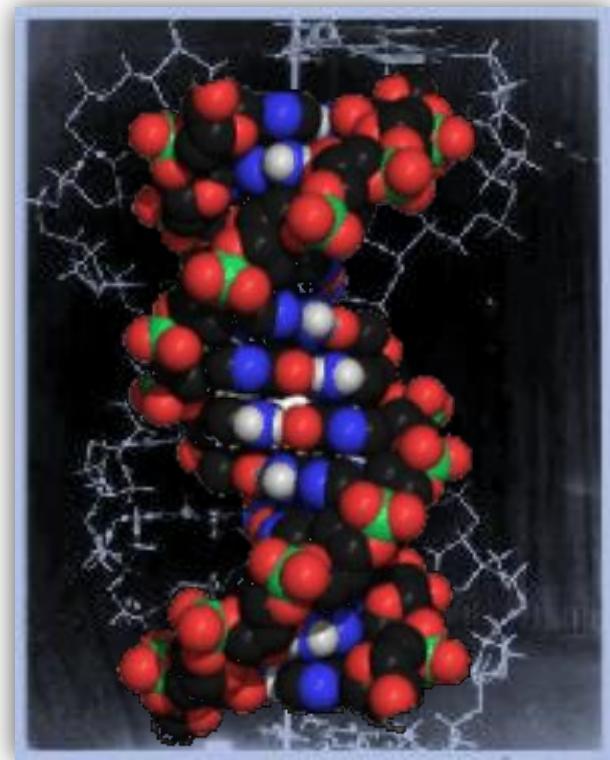




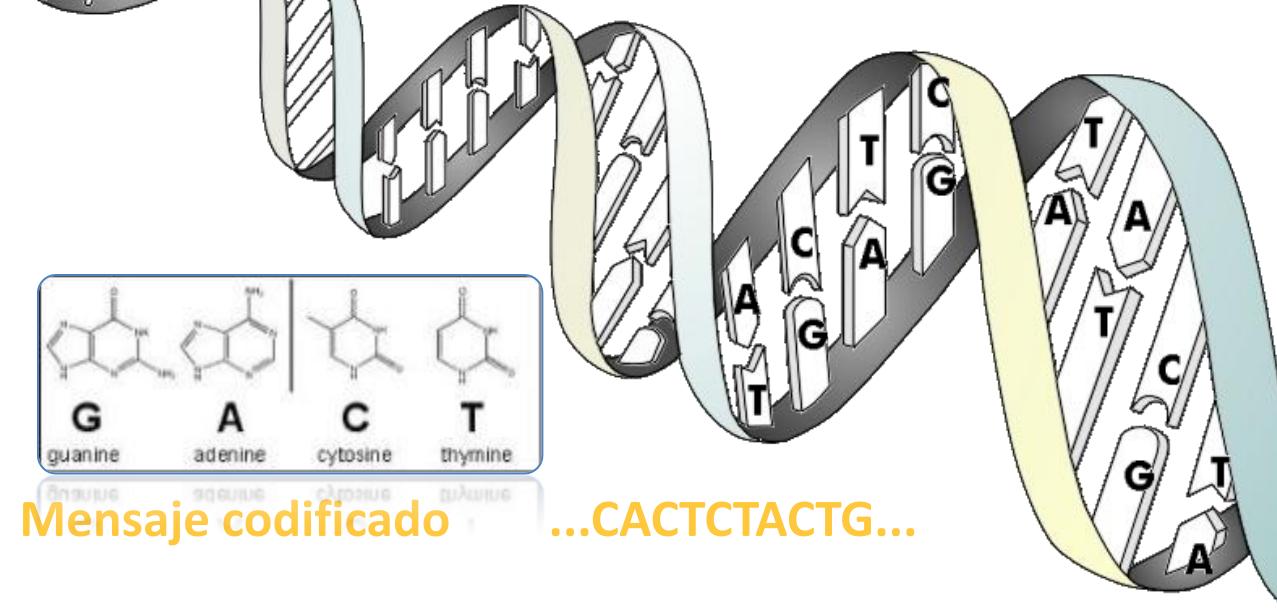
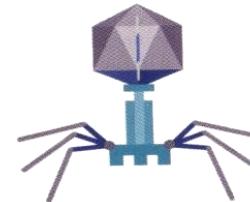
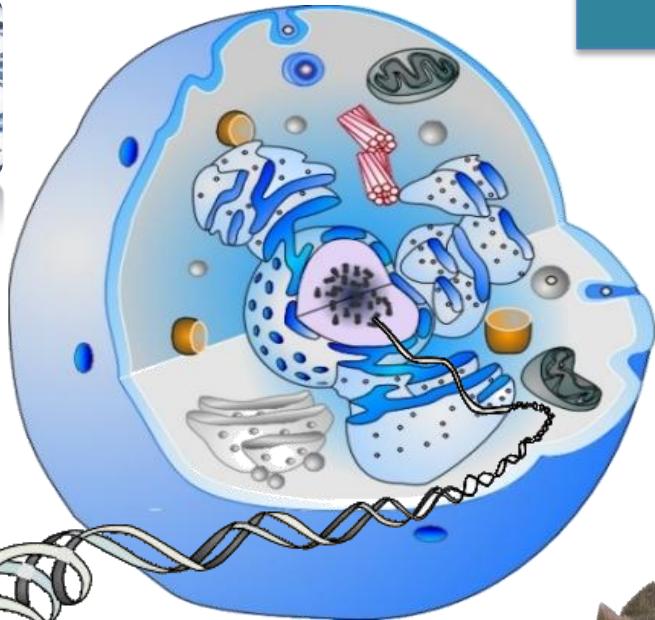
# DNA Doble Helix: the secret of life



1953

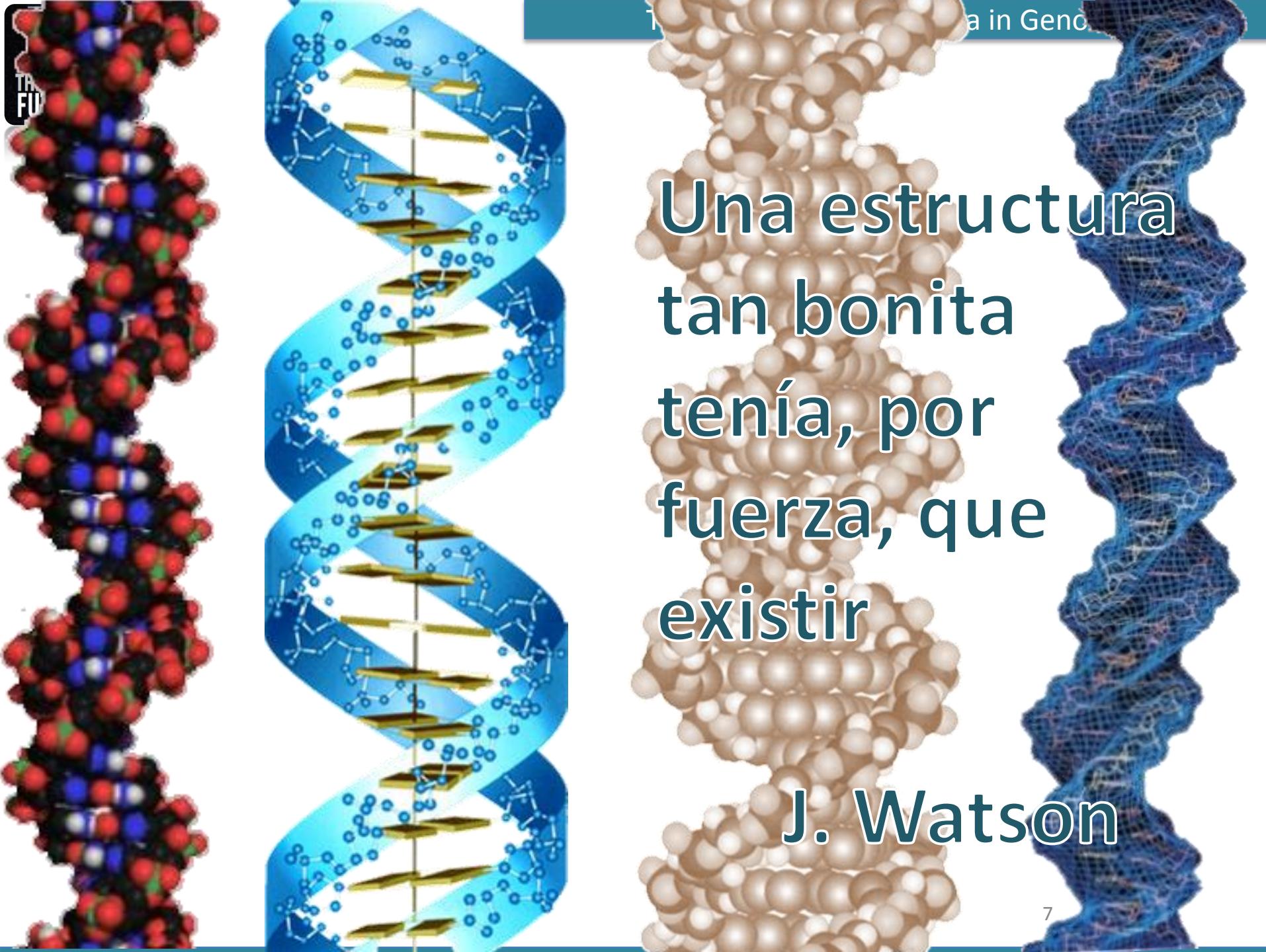


# The challenges of big data in Genomics



Mensaje codificado

...CACTCTACTG...



Tu herencia genética es la base de tu vida. La genética te da la información para crecer y desarrollarte. La genética te da la información para crecer y desarrollarte. La genética te da la información para crecer y desarrollarte.

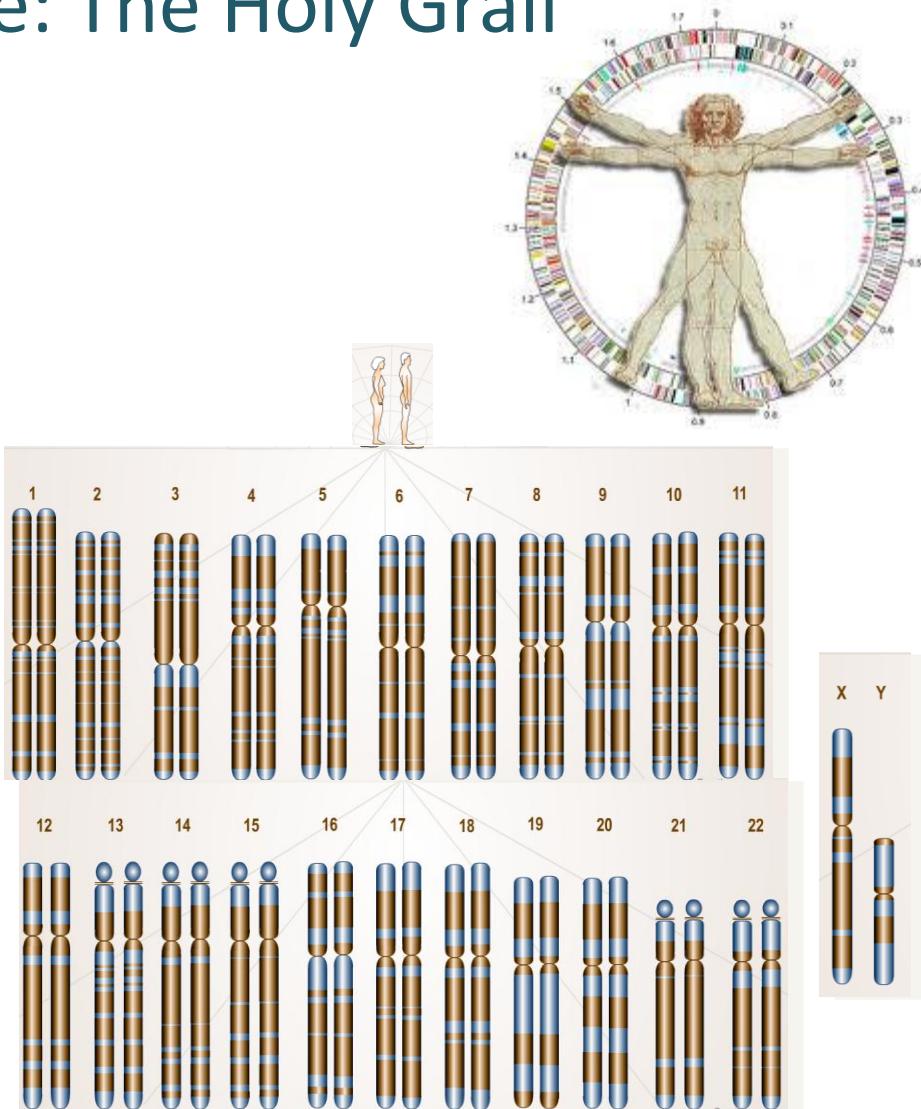
Una estructura  
tan bonita  
tenía, por  
fuerza, que  
existir

J. Watson

# The human genome: The Holy Grail



**DEFINICIÓN** Genoma humano  
**LOCUS** HSPG010101 3.158 Gb **ADN**  
**FECHA** 04-25-03  
**VERSIÓN** Ensamblaje 1.0  
**ORGANISMO** «Homo sapiens»  
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
 Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo.  
  
**TÍTULO** La secuencia  
**/fuente** 1, 3150000000  
**/cromosoma** "1-22, X, Y"  
**/nota** «Libro de la vida, santo grial, mapa humano»  
  
**INICIO**  
 1 agtcgcgtga gacttctgg accccgcacc aggctgtggg gtttctcaga taactgggcc  
 61 cctcgctca ggaggccctc accctctgt ctggtaaag ttcatggaa cagaaaaaaa  
 121 tggattatc tgcttcctcg gttaaagaag tacaaaaatg cattaatgcg atgcggaaaa  
 181 tcttagtgc tcccatctgt ctggatgtg tcaaggaaacc tgctccaca aagtgtgacc  
 241 acatatttg caaattttgc atgttgcggatc ttctcaaca gaagaaggcc ctttcacagt  
 301 gtcctttatg taagaatgtg ataaccaaaa ggaggccatca agaaagtacg agattttagtc  
 361 aacttgttgc agagcttattt aaaaatcattt gtgttttca gtttgacaca ggtttggagt  
 421 atgcaaaacag ctataatttt gcaaaaaaagg aaaaataactc tgcttgacat ctaaaagatg  
 481 aagtttctat catccaaatgt atggggctaca gaaaccgtgc caaagactt ctagacatgt  
 541 aaccggaaaa tccttccttg cagggaaatc gtctcgtgtt cttttttttt aaccttggaa  
 601 ctgtgagaac ttcggagaca aaggcagccgaa tacaactca aaggacttgc gtcacattt  
 661 aattgggatc tgattttctt gaagataccg ttaataaggc aacttattgc agtggggag  
 721 atcaagaattt gttacaaaatc accccctcaag gaaccaggga tgaaatcagt ttggattctg  
 781 caaaaaaaaaa gtcgtgttgc tttttttgcggatgttac aataactgaa catccatcaac  
 841 ccgttaataa ttgttgcggat accactgtgaa agcgtgcgc tgagggatc ccagaaaaatg  
 901 atcagggtatc ttctgttgc aacttgcgt tgggccatc tgccacaaaatcactatcc  
 961 gtcatttaca gcatggaaac agcagggttat tactactaa agacagaaatg aatgtaaaaaa  
 1021 aggtgttgcatt ctgttaataaa agcaaaacagc ctggcttgc aaggaggcca cataacagat  
 1081 gggctgttgc taaggaaaaa tggatgtataa ggccggactcc cagcacagaa aaaaggatg  
 1141 atctgttgc ttgtccccctg tggatgtggaa aagaatggaa taaggacaaaatcgt  
 1201 cagggatcatt tagatgtacta ggatgttgc ctggatatac aataatgcg agatccatcga  
 1261 aagttaatgtg tggttttcc agaaatgtg aacttgttgc ttctgtatc tcacatgtatc  
 1321 gggatgttgc atcaatgttgc aagaatgttgc atgttgcggatc cttttttttt  
 1381 aatattctgg ttcttcggat aaaaatgtact tactggccatc tgatccatcat  
 1441 tatgtaaaatg tggaaatgttgc cacttccaaat cttttttttt  
 1501 ttggggaaaaatc ctatccggaa agggccatcc tcccaactt aaggccatgt  
 1561 taattataggc agcatgttgc actggggccatcc agataatgcg agagcgtccc  
 1621 aattaaatggc taaaaggatc ctttccatgt gccttccatcc tgaggatattt  
 1681 cagattttgc agttccaaatg acttgcggatc tgatccatcat  
 1741 agaatgttgc agtgtatgttgc attactaata tggttgc  
 1801 ctattccatc tgaaaaatgttgc ctttccatcat  
 1861 aaacggaaaaatc tgaaatgttgc agggccatcc  
 1921 aacatccatcatc tgaaatgttgc agggccatcc  
 1981 atggccatccatc tgatccatcat  
 2041 ttgtatgttgc ttctgtatc  
 2101 tgatgttgc ttctgtatc  
 2161 aatgttgc ttctgtatc  
 2221 aatgttgc ttctgtatc  
 2281 aatgttgc ttctgtatc  
 2341 aatgttgc ttctgtatc  
 2401 aatgttgc ttctgtatc  
 2461 aatgttgc ttctgtatc  
 2521 aatgttgc ttctgtatc  
 2581 aatgttgc ttctgtatc  
 2641 aatgttgc ttctgtatc  
 2701 aatgttgc ttctgtatc  
 2761 aatgttgc ttctgtatc  
 2821 aatgttgc ttctgtatc  
 2881 aatgttgc ttctgtatc  
 2941 aatgttgc ttctgtatc  
 3001 aatgttgc ttctgtatc  
 3061 aatgttgc ttctgtatc  
 3121 aatgttgc ttctgtatc  
 3181 aatgttgc ttctgtatc  
 3241 aatgttgc ttctgtatc  
 3301 aatgttgc ttctgtatc  
 3361 aatgttgc ttctgtatc  
 3421 aatgttgc ttctgtatc  
 3481 aatgttgc ttctgtatc  
 3541 aatgttgc ttctgtatc  
 3601 aatgttgc ttctgtatc  
 3661 aatgttgc ttctgtatc  
 3721 aatgttgc ttctgtatc  
 3781 aatgttgc ttctgtatc  
 3841 aatgttgc ttctgtatc  
 3901 aatgttgc ttctgtatc  
 3961 aatgttgc ttctgtatc  
 4021 aatgttgc ttctgtatc  
 4081 aatgttgc ttctgtatc  
 4141 aatgttgc ttctgtatc  
 4201 aatgttgc ttctgtatc  
 4261 aatgttgc ttctgtatc  
 4321 aatgttgc ttctgtatc  
 4381 aatgttgc ttctgtatc  
 4441 aatgttgc ttctgtatc  
 4501 aatgttgc ttctgtatc  
 4561 aatgttgc ttctgtatc  
 4621 aatgttgc ttctgtatc  
 4681 aatgttgc ttctgtatc  
 4741 aatgttgc ttctgtatc  
 4801 aatgttgc ttctgtatc  
 4861 aatgttgc ttctgtatc  
 4921 aatgttgc ttctgtatc  
 4981 aatgttgc ttctgtatc  
 5041 aatgttgc ttctgtatc  
 5101 aatgttgc ttctgtatc  
 5161 aatgttgc ttctgtatc  
 5221 aatgttgc ttctgtatc  
 5281 aatgttgc ttctgtatc  
 5341 aatgttgc ttctgtatc  
 5401 aatgttgc ttctgtatc  
 5461 aatgttgc ttctgtatc  
 5521 aatgttgc ttctgtatc  
 5581 aatgttgc ttctgtatc  
 5641 aatgttgc ttctgtatc  
 5701 aatgttgc ttctgtatc  
 5761 aatgttgc ttctgtatc  
 5821 aatgttgc ttctgtatc  
 5881 aatgttgc ttctgtatc  
 5941 aatgttgc ttctgtatc  
 6001 aatgttgc ttctgtatc  
 6061 aatgttgc ttctgtatc  
 6121 aatgttgc ttctgtatc  
 6181 aatgttgc ttctgtatc  
 6241 aatgttgc ttctgtatc  
 6301 aatgttgc ttctgtatc  
 6361 aatgttgc ttctgtatc  
 6421 aatgttgc ttctgtatc  
 6481 aatgttgc ttctgtatc  
 6541 aatgttgc ttctgtatc  
 6601 aatgttgc ttctgtatc  
 6661 aatgttgc ttctgtatc  
 6721 aatgttgc ttctgtatc  
 6781 aatgttgc ttctgtatc  
 6841 aatgttgc ttctgtatc  
 6901 aatgttgc ttctgtatc  
 6961 aatgttgc ttctgtatc  
 7021 aatgttgc ttctgtatc  
 7081 aatgttgc ttctgtatc  
 7141 aatgttgc ttctgtatc  
 7201 aatgttgc ttctgtatc  
 7261 aatgttgc ttctgtatc  
 7321 aatgttgc ttctgtatc  
 7381 aatgttgc ttctgtatc  
 7441 aatgttgc ttctgtatc  
 7501 aatgttgc ttctgtatc  
 7561 aatgttgc ttctgtatc  
 7621 aatgttgc ttctgtatc  
 7681 aatgttgc ttctgtatc  
 7741 aatgttgc ttctgtatc  
 7801 aatgttgc ttctgtatc  
 7861 aatgttgc ttctgtatc  
 7921 aatgttgc ttctgtatc  
 7981 aatgttgc ttctgtatc  
 8041 aatgttgc ttctgtatc  
 8101 aatgttgc ttctgtatc  
 8161 aatgttgc ttctgtatc  
 8221 aatgttgc ttctgtatc  
 8281 aatgttgc ttctgtatc  
 8341 aatgttgc ttctgtatc  
 8401 aatgttgc ttctgtatc  
 8461 aatgttgc ttctgtatc  
 8521 aatgttgc ttctgtatc  
 8581 aatgttgc ttctgtatc  
 8641 aatgttgc ttctgtatc  
 8701 aatgttgc ttctgtatc  
 8761 aatgttgc ttctgtatc  
 8821 aatgttgc ttctgtatc  
 8881 aatgttgc ttctgtatc  
 8941 aatgttgc ttctgtatc  
 9001 aatgttgc ttctgtatc  
 9061 aatgttgc ttctgtatc  
 9121 aatgttgc ttctgtatc  
 9181 aatgttgc ttctgtatc  
 9241 aatgttgc ttctgtatc  
 9301 aatgttgc ttctgtatc  
 9361 aatgttgc ttctgtatc  
 9421 aatgttgc ttctgtatc  
 9481 aatgttgc ttctgtatc  
 9541 aatgttgc ttctgtatc  
 9601 aatgttgc ttctgtatc  
 9661 aatgttgc ttctgtatc  
 9721 aatgttgc ttctgtatc  
 9781 aatgttgc ttctgtatc  
 9841 aatgttgc ttctgtatc  
 9901 aatgttgc ttctgtatc  
 9961 aatgttgc ttctgtatc  
 10021 aatgttgc ttctgtatc  
 10081 aatgttgc ttctgtatc  
 10141 aatgttgc ttctgtatc  
 10201 aatgttgc ttctgtatc  
 10261 aatgttgc ttctgtatc  
 10321 aatgttgc ttctgtatc  
 10381 aatgttgc ttctgtatc  
 10441 aatgttgc ttctgtatc  
 10501 aatgttgc ttctgtatc  
 10561 aatgttgc ttctgtatc  
 10621 aatgttgc ttctgtatc  
 10681 aatgttgc ttctgtatc  
 10741 aatgttgc ttctgtatc  
 10801 aatgttgc ttctgtatc  
 10861 aatgttgc ttctgtatc  
 10921 aatgttgc ttctgtatc  
 10981 aatgttgc ttctgtatc  
 11041 aatgttgc ttctgtatc  
 11101 aatgttgc ttctgtatc  
 11161 aatgttgc ttctgtatc  
 11221 aatgttgc ttctgtatc  
 11281 aatgttgc ttctgtatc  
 11341 aatgttgc ttctgtatc  
 11401 aatgttgc ttctgtatc  
 11461 aatgttgc ttctgtatc  
 11521 aatgttgc ttctgtatc  
 11581 aatgttgc ttctgtatc  
 11641 aatgttgc ttctgtatc  
 11701 aatgttgc ttctgtatc  
 11761 aatgttgc ttctgtatc  
 11821 aatgttgc ttctgtatc  
 11881 aatgttgc ttctgtatc  
 11941 aatgttgc ttctgtatc  
 12001 aatgttgc ttctgtatc  
 12061 aatgttgc ttctgtatc  
 12121 aatgttgc ttctgtatc  
 12181 aatgttgc ttctgtatc  
 12241 aatgttgc ttctgtatc  
 12301 aatgttgc ttctgtatc  
 12361 aatgttgc ttctgtatc  
 12421 aatgttgc ttctgtatc  
 12481 aatgttgc ttctgtatc  
 12541 aatgttgc ttctgtatc  
 12601 aatgttgc ttctgtatc  
 12661 aatgttgc ttctgtatc  
 12721 aatgttgc ttctgtatc  
 12781 aatgttgc ttctgtatc  
 12841 aatgttgc ttctgtatc  
 12901 aatgttgc ttctgtatc  
 12961 aatgttgc ttctgtatc  
 13021 aatgttgc ttctgtatc  
 13081 aatgttgc ttctgtatc  
 13141 aatgttgc ttctgtatc  
 13201 aatgttgc ttctgtatc  
 13261 aatgttgc ttctgtatc  
 13321 aatgttgc ttctgtatc  
 13381 aatgttgc ttctgtatc  
 13441 aatgttgc ttctgtatc  
 13501 aatgttgc ttctgtatc  
 13561 aatgttgc ttctgtatc  
 13621 aatgttgc ttctgtatc  
 13681 aatgttgc ttctgtatc  
 13741 aatgttgc ttctgtatc  
 13801 aatgttgc ttctgtatc  
 13861 aatgttgc ttctgtatc  
 13921 aatgttgc ttctgtatc  
 13981 aatgttgc ttctgtatc  
 14041 aatgttgc ttctgtatc  
 14101 aatgttgc ttctgtatc  
 14161 aatgttgc ttctgtatc  
 14221 aatgttgc ttctgtatc  
 14281 aatgttgc ttctgtatc  
 14341 aatgttgc ttctgtatc  
 14401 aatgttgc ttctgtatc  
 14461 aatgttgc ttctgtatc  
 14521 aatgttgc ttctgtatc  
 14581 aatgttgc ttctgtatc  
 14641 aatgttgc ttctgtatc  
 14701 aatgttgc ttctgtatc  
 14761 aatgttgc ttctgtatc  
 14821 aatgttgc ttctgtatc  
 14881 aatgttgc ttctgtatc  
 14941 aatgttgc ttctgtatc  
 15001 aatgttgc ttctgtatc  
 15061 aatgttgc ttctgtatc  
 15121 aatgttgc ttctgtatc  
 15181 aatgttgc ttctgtatc  
 15241 aatgttgc ttctgtatc  
 15301 aatgttgc ttctgtatc  
 15361 aatgttgc ttctgtatc  
 15421 aatgttgc ttctgtatc  
 15481 aatgttgc ttctgtatc  
 15541 aatgttgc ttctgtatc  
 15601 aatgttgc ttctgtatc  
 15661 aatgttgc ttctgtatc  
 15721 aatgttgc ttctgtatc  
 15781 aatgttgc ttctgtatc  
 15841 aatgttgc ttctgtatc  
 15901 aatgttgc ttctgtatc  
 15961 aatgttgc ttctgtatc  
 16021 aatgttgc ttctgtatc  
 16081 aatgttgc ttctgtatc  
 16141 aatgttgc ttctgtatc  
 16201 aatgttgc ttctgtatc  
 16261 aatgttgc ttctgtatc  
 16321 aatgttgc ttctgtatc  
 16381 aatgttgc ttctgtatc  
 16441 aatgttgc ttctgtatc  
 16501 aatgttgc ttctgtatc  
 16561 aatgttgc ttctgtatc  
 16621 aatgttgc ttctgtatc  
 16681 aatgttgc ttctgtatc  
 16741 aatgttgc ttctgtatc  
 16801 aatgttgc ttctgtatc  
 16861 aatgttgc ttctgtatc  
 16921 aatgttgc ttctgtatc  
 16981 aatgttgc ttctgtatc  
 17041 aatgttgc ttctgtatc  
 17101 aatgttgc ttctgtatc  
 17161 aatgttgc ttctgtatc  
 17221 aatgttgc ttctgtatc  
 17281 aatgttgc ttctgtatc  
 17341 aatgttgc ttctgtatc  
 17401 aatgttgc ttctgtatc  
 17461 aatgttgc ttctgtatc  
 17521 aatgttgc ttctgtatc  
 17581 aatgttgc ttctgtatc  
 17641 aatgttgc ttctgtatc  
 17701 aatgttgc ttctgtatc  
 17761 aatgttgc ttctgtatc  
 17821 aatgttgc ttctgtatc  
 17881 aatgttgc ttctgtatc  
 17941 aatgttgc ttctgtatc  
 18001 aatgttgc ttctgtatc  
 18061 aatgttgc ttctgtatc  
 18121 aatgttgc ttctgtatc  
 18181 aatgttgc ttctgtatc  
 18241 aatgttgc ttctgtatc  
 18301 aatgttgc ttctgtatc  
 18361 aatgttgc ttctgtatc  
 18421 aatgttgc ttctgtatc  
 18481 aatgttgc ttctgtatc  
 18541 aatgttgc ttctgtatc  
 18601 aatgttgc ttctgtatc  
 18661 aatgttgc ttctgtatc  
 18721 aatgttgc ttctgtatc  
 18781 aatgttgc ttctgtatc  
 18841 aatgttgc ttctgtatc  
 18901 aatgttgc ttctgtatc  
 18961 aatgttgc ttctgtatc  
 19021 aatgttgc ttctgtatc  
 19081 aatgttgc ttctgtatc  
 19141 aatgttgc ttctgtatc  
 19201 aatgttgc ttctgtatc  
 19261 aatgttgc ttctgtatc  
 19321 aatgttgc ttctgtatc  
 19381 aatgttgc ttctgtatc  
 19441 aatgttgc ttctgtatc  
 19501 aatgttgc ttctgtatc  
 19561 aatgttgc ttctgtatc  
 19621 aatgttgc ttctgtatc  
 19681 aatgttgc ttctgtatc  
 19741 aatgttgc ttctgtatc  
 19801 aatgttgc ttctgtatc  
 19861 aatgttgc ttctgtatc  
 19921 aatgttgc ttctgtatc  
 19981 aatgttgc ttctgtatc  
 20041 aatgttgc ttctgtatc  
 20101 aatgttgc ttctgtatc  
 20161 aatgttgc ttctgtatc  
 20221 aatgttgc ttctgtatc  
 20281 aatgttgc ttctgtatc  
 20341 aatgttgc ttctgtatc  
 20401 aatgttgc ttctgtatc  
 20461 aatgttgc ttctgtatc  
 20521 aatgttgc ttctgtatc  
 20581 aatgttgc ttctgtatc  
 20641 aatgttgc ttctgtatc  
 20701 aatgttgc ttctgtatc  
 20761 aatgttgc ttctgtatc  
 20821 aatgttgc ttctgtatc  
 20881 aatgttgc ttctgtatc  
 20941 aatgttgc ttctgtatc  
 21001 aatgttgc ttctgtatc  
 21061 aatgttgc ttctgtatc  
 21121 aatgttgc ttctgtatc  
 21181 aatgttgc ttctgtatc  
 21241 aatgttgc ttctgtatc  
 21301 aatgttgc ttctgtatc  
 21361 aatgttgc ttctgtatc  
 21421 aatgttgc ttctgtatc  
 21481 aatgttgc ttctgtatc  
 21541 aatgttgc ttctgtatc  
 21601 aatgttgc ttctgtatc  
 21661 aatgttgc ttctgtatc  
 21721 aatgttgc ttctgtatc  
 21781 aatgttgc ttctgtatc  
 21841 aatgttgc ttctgtatc  
 21901 aatgttgc ttctgtatc  
 21961 aatgttgc ttctgtatc  
 22021 aatgttgc ttctgtatc  
 22081 aatgttgc ttctgtatc  
 22141 aatgttgc ttctgtatc  
 22201 aatgttgc ttctgtatc  
 22261 aatgttgc ttctgtatc  
 22321 aatgttgc ttctgtatc  
 22381 aatgttgc ttctgtatc  
 22441 aatgttgc ttctgtatc  
 22501 aatgttgc ttctgtatc  
 22561 aatgttgc ttctgtatc  
 22621 aatgttgc ttctgtatc  
 22681 aatgttgc ttctgtatc  
 22741 aatgttgc ttctgtatc  
 22801 aatgttgc ttctgtatc  
 22861 aatgttgc ttctgtatc  
 22921 aatgttgc ttctgtatc  
 22981 aatgttgc ttctgtatc  
 23041 aatgttgc ttctgtatc  
 23101 aatgttgc ttctgtatc  
 23161 aatgttgc ttctgtatc  
 23221 aatgttgc ttctgtatc  
 23281 aatgttgc ttctgtatc  
 23341 aatgttgc ttctgtatc  
 23401 aatgttgc ttctgtatc  
 23461 aatgttgc ttctgtatc  
 23521 aatgttgc ttctgtatc  
 23581 aatgttgc ttctgtatc  
 23641 aatgttgc ttctgtatc  
 23701 aatgttgc ttctgtatc  
 23761 aatgttgc ttctgtatc  
 23821 aatgttgc ttctgtatc  
 23881 aatgttgc ttctgtatc  
 23941 aatgttgc ttctgtatc  
 24001 aatgttgc ttctgtatc  
 24061 aatgttgc ttctgtatc  
 24121 aatgttgc ttctgtatc  
 24181 aatgttgc ttctgtatc  
 24241 aatgttgc ttctgtatc  
 24301 aatgttgc ttctgtatc  
 24361 aatgttgc ttctgtatc  
 24421 aatgttgc ttctgtatc  
 24481 aatgttgc ttctgtatc  
 24541 aatgttgc ttctgtatc  
 24601 aatgttgc ttctgtatc  
 24661 aatgttgc ttctgtatc  
 24721 aatgttgc ttctgtatc  
 24781 aatgttgc ttctgtatc  
 24841 aatgttgc ttctgtatc  
 24901 aatgttgc ttctgtatc  
 24961 aatgttgc ttctgtatc  
 25021 aatgttgc ttctgtatc  
 25081 aatgttgc ttctgtatc  
 25141 aatgttgc ttctgtatc  
 25201 aatgttgc ttctgtatc  
 25261 aatgttgc ttctgtatc  
 25321 aatgttgc ttctgtatc  
 25381 aatgttgc ttctgtatc  
 25441 aatgttgc ttctgtatc  
 25501 aatgttgc ttctgtatc  
 25561 aatgttgc ttctgtatc  
 25621 aatgttgc ttctgtatc  
 25681 aatgttgc ttctgtatc  
 25741 aatgttgc ttctgtatc  
 25801 aatgttgc ttctgtatc  
 25861 aatgttgc ttctgtatc  
 25921 aatgttgc ttctgtatc  
 25981 aatgttgc ttctgtatc  
 26041 aatgttgc ttctgtatc  
 26101 aatgttgc ttctgtatc  
 26161 aatgttgc ttctgtatc  
 26221 aatgttgc ttctgtatc  
 26281 aatgttgc ttctgtatc  
 26341 aatgttgc ttctgtatc  
 26401 aatgttgc ttctgtatc  
 26461 aatgttgc ttctgtatc  
 26521 aatgttgc ttctgtatc  
 26581 aatgttgc ttctgtatc  
 26641 aatgttgc ttctgtatc  
 26701 aatgttgc ttctgtatc  
 26761 aatgttgc ttctgtatc  
 26821 aatgttgc ttctgtatc  
 26881 aatgttgc ttctgtatc  
 26941 aatgttgc ttctgtatc  
 27001 aatgttgc ttctgtatc  
 27061 aatgttgc ttctgtatc  
 27121 aatgttgc ttctgtatc  
 27181 aatgttgc ttctgtatc  
 27241 aatgttgc ttctgtatc  
 27301 aatgttgc ttctgtatc  
 27361 aatgttgc ttctgtatc  
 27421 aatgttgc ttctgtatc  
 27481 aatgttgc ttctgtatc  
 27541 aatgttgc ttctgtatc  
 27601 aatgttgc ttctgtatc  
 27661 aatgttgc ttctgtatc  
 27721 aatgttgc ttctgtatc  
 27781 aatgttgc ttctgtatc  
 27841 aatgttgc ttctgtatc  
 27901 aatgttgc ttctgtatc  
 27961 aatgttgc ttctgtatc  
 28021 aatgttgc ttctgtatc  
 28081 aatgttgc ttctgtatc  
 28141 aatgttgc ttctgtatc  
 28201 aatgttgc ttctgtatc  
 28261 aatgttgc ttctgtatc  
 28321 aatgttgc ttctgtatc  
 28381 aatgttgc ttctgtatc  
 28441 aatgttgc ttctgtatc  
 28501 aatgttgc ttctgtatc  
 28561 aatgttgc ttctgtatc  
 28621 aatgttgc ttctgtatc  
 28681 aatgttgc ttctgtatc  
 28741 aatgttgc ttctgtatc  
 28801 aatgttgc ttctgtatc  
 28861 aatgttgc ttctgtatc  
 28921 aatgttgc ttctgtatc  
 28981 aatgttgc ttctgtatc  
 29041 aatgttgc ttctgtatc  
 29101 aatgttgc ttctgtatc  
 29161 aatgttgc ttctgtatc  
 29221 aatgttgc ttctgtatc  
 29281 aatgttgc ttctgtatc  
 29341 aatgttgc ttctgtatc  
 29401 aatgttgc ttctgtatc  
 29461 aatgttgc ttctgtatc  
 29521 aatgttgc ttctgtatc  
 29581 aatgttgc ttctgtatc  
 29641 aatgttgc ttctgtatc  
 29701 aatgttgc ttctgtatc  
 29761 aatgttgc ttctgtatc  
 29821 aatgttgc ttctgtatc  
 29881 aatgttgc ttctgtatc  
 29941 aatgttgc ttctgtatc  
 30001 aatgttgc ttctgtatc  
 30061 aatgttgc ttctgtatc  
 30121 aatgttgc ttctgtatc  
 30181 aatgttgc ttctgtatc  
 30241 aatgttgc ttctgtatc  
 30301 aatgttgc ttctgtatc  
 30361 aatgttgc ttctgtatc  
 30421 aatgttgc ttctgtatc  
 30481 aatgttgc ttctgtatc  
 30541 aatgttgc ttctgtatc  
 30601 aatgttgc ttctgtatc  
 30661 aatgttgc ttctgtatc  
 30721 aatgttgc ttctgtatc  
 30781 aatgttgc ttctgtatc  
 30841 aatgttgc ttctgtatc  
 30901 aatgttgc ttctgtatc  
 30961 aatgttgc ttctgtatc  
 31021 aatgttgc ttctgtatc  
 31081 aatgttgc ttctgtatc  
 31141 aatgttgc ttctgtatc  
 31201 aatgttgc ttctgtatc  
 31261 aatgttgc ttctgtatc  
 31321 aatgttgc ttctgtatc  
 31381 aatgttgc ttctgtatc  
 31441 aatgttgc ttctgtatc  
 31501 aatgttgc ttctgtatc  
 31561 aatgttgc ttctgtatc  
 31621 aatgttgc ttctgtatc  
 31681 aatgttgc ttctgtatc  
 31741 aatgttgc ttctgtatc  
 31801 aatgttgc ttctgtatc  
 31861 aatgttgc ttctgtatc  
 31921 aatgttgc ttctgtatc  
 31981 aatgttgc ttctgtatc  
 32041 aatgttgc ttctgtatc  
 32101 aatgttgc ttctgtatc  
 32161 aatgttgc ttctgtatc  
 32221 aatgttgc ttctgtatc  
 32281 aatgttgc ttctgtatc  
 32341 aatgttgc ttctgtatc  
 32401 aatgttgc ttctgtatc  
 32461 aatgttgc ttctgtatc  
 32521 aatgttgc ttctgtatc  
 32581 aatgttgc ttctgtatc  
 32641 aatgttgc ttctgtatc  
 32701 aatgttgc ttctgtatc  
 32761 aatgttgc ttctgtatc  
 32821 aatgttgc ttctgtatc  
 32881 aatgttgc ttctgtatc  
 32941 aatgttgc ttctgtatc  
 33001 aatgttgc ttctgtatc  
 33061 aatgttgc ttctgtatc  
 33121 aatgttgc ttctgtatc  
 33181 aatgttgc ttctgtatc  
 33241 aatgttgc ttctgtatc  
 33301 aatgttgc ttctgtatc  
 33361 aatgttgc ttctgtatc  
 33421 aatgttgc ttctgtatc  
 33481 aatgttgc ttctgtatc  
 33541 aatgttgc ttctgtatc  
 33601 aatgttgc ttctgtatc  
 33661 aatgttgc ttctgtatc  
 33721 aatgttgc ttctgtatc  
 33781 aatgttgc ttctgtatc  
 33841 aatgttgc ttctgtatc  
 33901 aatgttgc ttctgtatc  
 33961 aatgttgc ttctgtatc  
 34021 aatgttgc ttctgtatc  
 34081 aatgttgc ttctgtatc  
 34141 aatgttgc ttctgtatc  
 34201 aatgttgc ttctgtatc  
 34261 aatgttgc ttctgtatc  
 34321 aatgttgc ttctgtatc  
 34381 aatgttgc ttctgtatc  
 34441 aatgttgc ttctgtatc  
 34501 aatgttgc ttctgtatc  
 34561 aatgttgc ttctgtatc  
 34621 aatgttgc ttctgtatc  
 34681 aatgttgc ttctgtatc  
 34741 aatgttgc ttctgtatc  
 34801 aatgttgc ttctgtatc  
 34861 aatgttgc ttctgtatc  
 34921 aatgttgc ttctgtatc  
 34981 aatgttgc ttctgtatc  
 35041 aatgttgc ttctgtatc  
 35101 aatgttgc ttctgtatc  
 35161 aatgttgc ttctgtatc  
 35221 aatgttgc ttctgtatc  
 35281 aatgttgc ttctgtatc  
 35341 aatgttgc ttctgtatc  
 35401 aatgttgc ttctgtatc  
 35461 aatgttgc ttctgtatc  
 35521 aatgttgc ttctgtatc  
 35581 aatgttgc ttctgtatc  
 35641 aatgttgc ttctgtatc  
 35701 aatgttgc ttctgtatc  
 35761 aatgttgc ttctgtatc  
 35821 aatgttgc ttctgtatc  
 35881 aatgttgc ttctgtatc  
 35941 aatgttgc ttctgtatc  
 36001 aatgttgc ttctgtatc  
 36061 aatgttgc ttctgtatc  
 36121 aatgttgc ttctgtatc  
 36181 aatgttgc ttctgtatc  
 36241 aatgttgc ttctgtatc  
 36301 aatgttgc ttctgtatc  
 36361 aatgttgc ttctgtatc  
 36421 aatgttgc ttctgtatc  
 36481 aatgttgc ttctgtatc  
 36541 aatgttgc ttctgtatc  
 36601 aatgttgc ttctgtatc  
 36661 aatgttgc ttctgtatc  
 36721 aatgttgc ttctgtatc  
 36781 aatgttgc ttctgtatc  
 36841 aatgttgc ttctgtatc  
 36901 aatgttgc ttctgtat

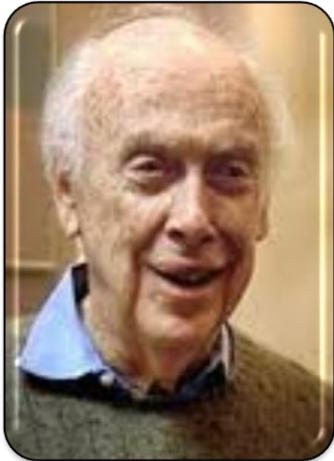


*The mapping of the human genome is ‘the greatest intellectual moment in history.’*



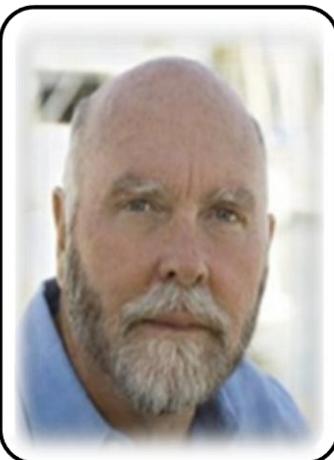
Matt Ridley

ACTGA  
CTTACG



*It's a giant resource that will change mankind, like the printing press*

James Watson



*This period is a very historic time, a new starting point*

Craig Venter

# The achievements of the HGP has radically changed the practice of biomedical research

## Distinguishing characteristics of Genomics

- Big teams -> Multidisciplinary and international teams
- New high-throughput technologies for large-scale data production (Omics)
- “Discovery science” or “Data driven” approach vs. “hypothesis driven” approach
- Computational intensity and expertise
- High standard for data quality
- Rapid data release
- Attention to societal implications



### International Human Genome Sequencing Consortium\*

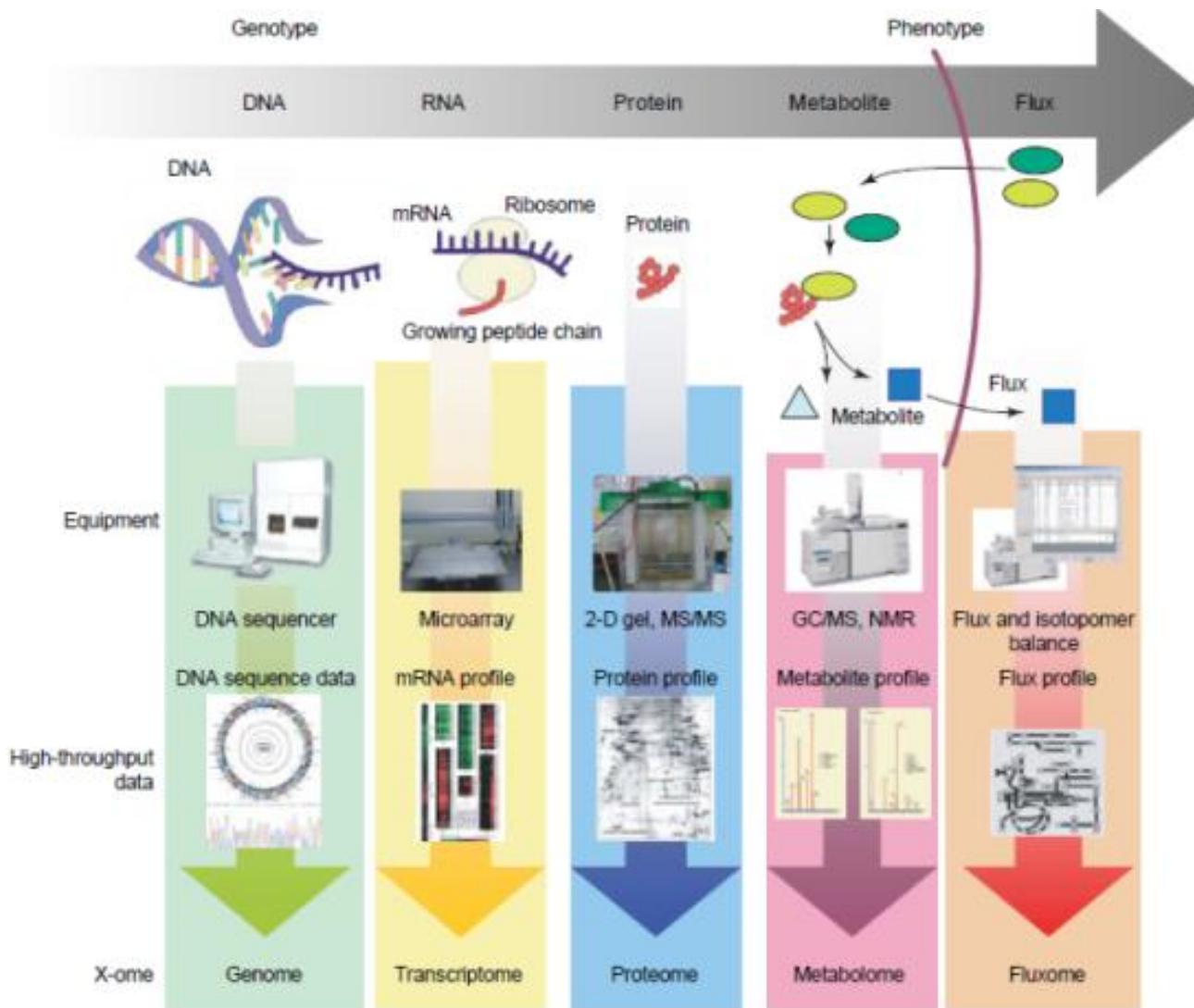
<b>Genome Sequencing Centers</b> (listed in order of total genomic sequence contributed, with a partial list of personnel. A full list of contributors at each centre is available as Supplementary Information.)
Wellcome Trust Sanger Institute: Eric V. Lander*, James M. Lemire*, Anne Kersey*, Christopher Clark*, David J. Zody*, Jennifer Johnson*, Karen E. Cross, Kim Doherty, Michael Drury*, William Platig*, Paul Pielzus*, Diane Rappaport*, Karen Reilly*, Andrew Rutherford*, John Auerbach*, Leah Karpel*, Jessica Lehoczky*, Alison Levine*, Paul McEwan*, Karin McMurrae*, James McMurrae*, JEB P. Mehta*, Ober Mirenski*, William Morris*, Jerome Rozen*, Andrew Shatkin*, Carla Souza*, Nicole Sturge-Thomason*, Nicole Stepaniak*, Aravind Subramanian*, & Bradley Wray*
The Sanger Center: Jane Rogers*, John Sulston*, Richard Attwells*, Stephen Beck*, David Bentley*, John Burton*, Christopher Clark*, Nigel Carter*, Alan Cawdron*, Rebecca Chapman*, Lucas Dekker*, Andrew Durbin*, Ian Dunham*, Richard Durbin*, Luis Freire*, Darren Grimes*, Steven Gregory*, Tim Hubbard*, Ben Humphries*, Adrienne Hunt*, Matthew Jones*, Christine Lucy*, Amanda McMurray*, Lucy Mcmurray*, Simon Morris*, Sarah Miller*, James G. Mullikin*, Andrew Mungall*, Robert Mungall*, Mark Rose*, Katrina Shawcross*, Jason Stuker*, & Sarah Storey*
Washington University Genome Sequencing Center: Robert K. Wilson*, Richard E. Wilson*, Lukasz W. Wilusz*, John S. McPherson*, Marcus A. Maruy*, Elizabet B. Morley*, Leandra A. Putzer*, Asaf T. Orenstein*, Kenneth R. Peay*, Warren R. Gilot*, Stephanek R. Chiszar*, Michael L. Venter*, Ken C. Deininger*, Tracey L. Mies*, Andrew Deininger*, Jason B. Kramer*, Lisa L. Cook*, Robert S. Fuller*, Douglas J. Johnson*, Patrick J. Hillis*, Sandra M. Clinton*
US DOE Joint Genome Institute: Trevor Harkness*, Scott Hannon*, Tom Hennell*, Warren Boppert*, Ann-Ping Cheng*, Anne Ober*, Susan Lucas*, Christopher Rohr*, Colleen Stepaniak*, & Nancy Prater*
Naylor College of Medicine Human Genome Sequencing Center: Michael A. Dray*, Dennis M. Murphy*, Steven S. Scherer*, John S. Aszkenasy*, Eric J. Soderberg*, Alan C. Umey*, Catherine M. River*, James R. Gersoff*, Michael L. Meckler*, Steven L. Taylor*, Riley S. KarcherLapin*, David L. Nelson, & George M. Weissenbach*
NIH/National Genome Sciences Center: Yoshiyuki Sakaki*, Asao Fujiyama*, Hisaharu Kubota*, Tetsuya Fukaz*, Atsushi Toyoda*, Toshihiko Mori*, Ohshiro Kawagoe*, Natsuo Matsuzaki*, Fumio Tobe* & Takuji Hayashi*
Sanger/CBBS JMR-BS2C: Jean Winterton*, Roland Holley*, William Saunier*, Francois Artigau*, Philippe Brattain*, Thomas Brabec*, Eric Pellegrini*, Catherine Hubert*, & Patrick Winsor*
NYC Sequencing Center: Douglas B. Smith*, Lynn Doucette-Stamm*, Marc Baldwin*, Keith Weintraub*, Hong Mai Lu*, & John Detrait*
Department of Genetics Analysis, Institute of Molecular
Genetics: Andris Rosenthal*, Matthias Pospisil*, Gerald Rybníkár*, Steffen Taubert*, & Andreas Rump*
Beijing Genomics Institute/Human Genome Center: Haixiang Yang*, Jun Yu*, Jun Wang*, Suyong Huang*, & Jun Gu*
Human Genome Sequencing Center, The Institute for Systems Biology: Leroy Hood*, Lee Ritenour*, Amos Madar*, & Shoukhrat Mirza*
Stanford Genome Technology Center: Ronald W. Davis*, Randy A. Fedarko*, A. Pa Abdol*, & Michael J. Petrus*
Stanford Human Genome Center: Richard M. Myers*, Jeremy Schatz*, Mark Dickson*, Jane Greenwood*, & Daniel R. Sod*
University of Washington Genome Center: Maynard W. Olson*, Rajendra Kad*, & Christopher Raymond*
Department of Molecular Biology, Yale University School of Medicine: Michaelis Schloss*, Kathleen Kennedy*, & Steven Mountsin*
University of Texas Southwestern Medical Center of Dallas: Uta E. Evans*, Maria Athanasiadis*, & Roger Schulz*
University of Oklahoma's Advanced Center for Genome Technology: Bruce A. Pfeifer*, Feng Chen*, & Hoang Phu*
Max Planck Institute for Molecular Genetics: Jürgen Richter*, Hans Lehrach*, & Ralford Pfeiffer*
GATC Spring Porter Laboratory, Uta Anneliese Hasan Genome Center: Michaela Röschenthal*, Melitta de la Riva*, & Helmut Bröcker*
CBP—German Research Centre for Biotechnology: Helmut Stölzel*, Klaus Herremans*, & Gabriele Kortes*
Genome Analytics Group (listed in alphabetical order, also includes individuals listed under heading): Richard Agarwala*, J. Aszkenasy*, Jeffrey A. Barker*, Kirk Baker*, Serafini Scagliuzzi*, David Birney*, Peter Bork*, Daniel S. Brown*, Christopher B. Burtt*, Lorenzo Cerruti*, Hilde-Christen Choi*, Debra Church*, Michael Claverie*, Richard P. Cooley*, Tolka Drorika*, Sean E. Eddy*, Evan E. Fischer*, Terrence S. Furey*, James Gogos*, James C. R. Gifford*, Guyana Hanmer*, Yoshitaka Hayashizaki*, David Housset*, Henning Hermjakob*, Kunihiro Hoshino*, Werner Jähn*, L. Steven Johnson*, Thomas A. Jones*, Simon Ross*, Ark Karpas*, Sean Kennedy*, R. Jones*, Paul Kitz*, Eugene K. Koonin*, Ian Korf*, Daniel Kugy*, David Lauter*, Todd M. Lewis*, Adam McElroy*, Teijo Mikkelsen*, John V. Moran*, Woods Müller*, Vicki J. Patrus*, Chris P. Ponting*, Greg Schuler*, Jing Shiu*, Guy Stalter*, Arjan A. Smits*, Elia Stupka*, Joseph Szostakowski*, Danielle Thivierge-Ming*, James Tsui*, May*, Lucas Wayner*, John Weller*, Raymond Wimmer*, Abby Yilmaz*, Paul L. Wolf*, Kenneth L. Wohlf*, Shou-Ying Yang*, & Fei-Ying Yeh*
Scientific management: National Human Genome Research Institute, US National Institutes of Health: Francis Collins*, Mark S. Sayers*, Jane Peterson*, Adam Pellegrini*, & Kris A. Wetterhahn*, Office of Science, US Department of Energy: Aristide Patrinos*, Wellcome Trust: Michael J. Hogan*



Barcelona Computing Center (BCS)



# Omics



# The achievements of the HGP has radically changed the practice of biomedical research

## Distinguishing characteristics of Genomics

### International Human Genome Sequencing Consortium\*

Genome Sequencing Centers (Listed in order of total genomic sequence contributed, with a partial list of personnel. A full list of contributors at each centre is available as Supplementary Information.)
Wellcome Trust Sanger Institute: Eric V. Lander <sup>1</sup> , James M. Lehoczky <sup>1</sup> , Bruce Birren <sup>1</sup> , Christopher Clark <sup>1</sup> , Zandy C. Collins <sup>1</sup> , Robert D. Edwards <sup>1</sup> , Kim Eustace <sup>1</sup> , Kim Fitch <sup>1</sup> , Michael Gruis <sup>1</sup> , William Hefner <sup>1</sup> , Kent Hockley <sup>1</sup> , Kim Jones <sup>1</sup> , Michael Karch <sup>1</sup> , William Platzer <sup>1</sup> , Paul Pukac <sup>1</sup> , Diane Rapp <sup>1</sup> , Karen Reilly <sup>1</sup> , Andrew Redmond <sup>1</sup> , John Rawlins <sup>1</sup> , Lisa Saad <sup>1</sup> , Jessica Schatz <sup>1</sup> , Robin Seeger <sup>1</sup> , Paul St. Onge <sup>1</sup> , David T. Tewari <sup>1</sup> , Michael Tipton <sup>1</sup> , John Walker <sup>1</sup> , Richard Wiltshire <sup>1</sup> , James Woodward <sup>1</sup> , JEB P. Young <sup>1</sup> , Ober Witzel <sup>1</sup> , Karen McCarren <sup>1</sup> , Jerome Rozen <sup>1</sup> , Christopher Schuster <sup>1</sup> , Carla Souza <sup>1</sup> , Nicole Stange-Thomassen <sup>1</sup> , Nicole Stepaniak <sup>1</sup> , Aravind Subramanian <sup>1</sup> & Bradley Wray <sup>1</sup>
University of Washington Genome Center: Raymond W. Durbin <sup>2</sup> , Rajender Kalsi <sup>2</sup> & Christopher Raymond <sup>2</sup>
Department of Molecular Biology, Yale University School of Medicine: Michaelay Shain <sup>3</sup> , Kathleen Kennedy <sup>3</sup> & Steven Mountz <sup>3</sup>
University of Texas Southwestern Medical Center at Dallas: Uta Franke <sup>4</sup> , Maria Attanasio <sup>4</sup> & Roger Schultz <sup>4</sup>
University of Oklahoma's Advanced Center for Genome Technology: Eric P. Feingold <sup>5</sup> , Greg Ober <sup>5</sup> & Kenneth Pohl <sup>5</sup>
Max Planck Institute for Molecular Genetics: Jürgen Richter <sup>6</sup> , Hans Lehrach <sup>6</sup> & Rüdiger von Knebel <sup>6</sup>
Cold Spring Harbor Laboratory, Lila Wallace Human Genome Center: William R. McCombie <sup>7</sup> , Melissa de la Torre <sup>7</sup> & Holley Drabkin <sup>7</sup>
CRG—Catalan Research Centre for Biotechnology: Helmut Rückert <sup>8</sup> , Massimiliano Riccardi <sup>8</sup> & Gabriele Novembre <sup>8</sup>
* Genome Analysis Group (Listed in alphabetical order, also includes individuals listed under other headings):
Rick Agarwala <sup>1</sup> , Ian Ainsworth <sup>1</sup> , Jeffrey A. Barker <sup>1</sup> , Kirk Baker <sup>1</sup> , Serafin Salzberg <sup>1</sup> , Evan Birney <sup>1</sup> , Peter Bork <sup>1</sup> , Daniel G. Brown <sup>1</sup> , Christopher B. Burtt <sup>1</sup> , Lorenzo Cerruti <sup>1</sup> , Hye-Jeann Choi <sup>1</sup> , Debra Church <sup>1</sup> , Michael Clancy <sup>1</sup> , Richard R. Coplin <sup>1</sup> , Tolka Drabkin <sup>1</sup> , Sean E. Eddy <sup>1</sup> , Evan C. Fischer <sup>1</sup> , Terence S. Fyffe <sup>1</sup> , James Gaggen <sup>1</sup> , James C. E. Gilberg <sup>1</sup> , Guy Hanafi <sup>1</sup> , Yoshinobu Hayashizaki <sup>1</sup> , David Housman <sup>1</sup> , Henning Hermjakob <sup>1</sup> , Karsten Hohmann <sup>1</sup> , Wimerae Jones <sup>1</sup> , L. Steven Johnson <sup>1</sup> , Thomas A. Jones <sup>1</sup> , Simon Raff <sup>1</sup> , Ann Kasprzyk <sup>1</sup> , Sean Kennedy <sup>1</sup> , R. Jones <sup>1</sup> , Paul Kitz <sup>1</sup> , Eugene K. Koon <sup>1</sup> , Ian Korf <sup>1</sup> , Daniel Kugy <sup>1</sup> , David Lauter <sup>1</sup> , Todd M. Lewis <sup>1</sup> , Adam McJilton <sup>1</sup> , Teijo Mikkelsen <sup>1</sup> , John V. Minas <sup>1</sup> , Woods Mukler <sup>1</sup> , Vicki J. Paterson <sup>1</sup> , Chris P. Ponting <sup>1</sup> , Greg Schuler <sup>1</sup> , Jing Shiu <sup>1</sup> , Guy Stalter <sup>1</sup> , Arun P. A. Smith <sup>1</sup> , Elia Stupka <sup>1</sup> , Joseph Szostakowski <sup>1</sup> , Danielle Thibault-Ming <sup>1</sup> , Jason Thompson-May <sup>1</sup> , Lucas Wayner <sup>1</sup> , John Wilkes <sup>1</sup> , Raymond Wilson <sup>1</sup> , Asya Yilmaz <sup>1</sup> , Paul J. White <sup>1</sup> , Kenneth R. Wolfe <sup>1</sup> , Shou-Ping Yang <sup>1</sup> & Ru-fang Yin <sup>1</sup>
Scientific management: National Human Genome Research Institute, US National Institutes of Health: Francis Collins <sup>9</sup> , Mark S. Sayers <sup>9</sup> , Jane Peterson <sup>9</sup> , Adam Pehlivanian <sup>9</sup> & Kris A. Wetterhahn <sup>9</sup> ; Office of Science, US Department of Energy: Aristide Patrinos <sup>10</sup> ; The Wellcome Trust: Michael J. Morgan <sup>10</sup>

- Big teams -> Multidisciplinary and international teams
- New high-throughput technologies for large-scale data production (Omics)
- “Discovery science” or “Data driven” approach vs. “hypothesis driven” approach
- Computational intensity and expertise
- High standard for data quality
- Rapid data release
- Attention to societal implications



Barcelona Computing Center (BCS)

# BIG DATA SCIENCE

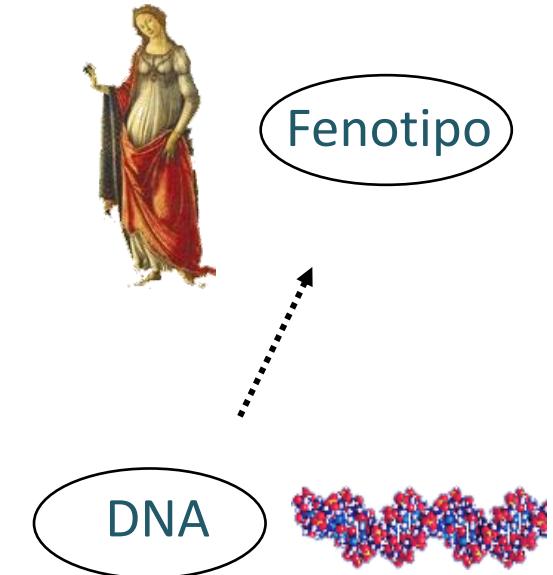


# The Fundamental Question



Clarividencia

René Magritte



¿Cómo se decodifica el  
mensaje genético para  
formar el fenotipo?



The HGP has changed the way we conceptualize molecular biology

## The HGP has had a profound consequence in the conceptualization of biological systems

New paradigm

Biology as a  
Informational Science

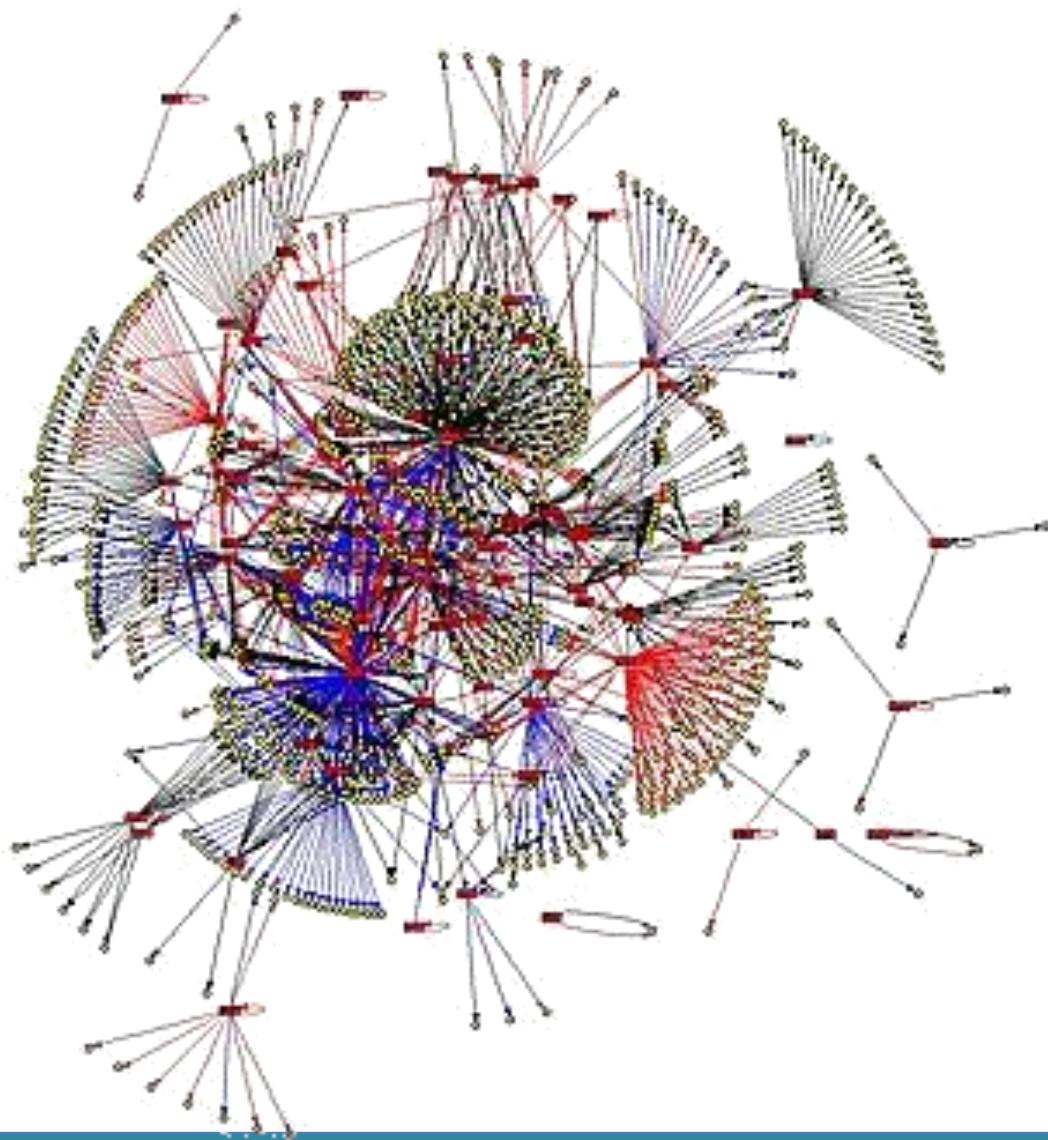


- The analysis of biological systems in terms of the storage, transmission and transformation of the information coded in the genomes



# Integrative Biology

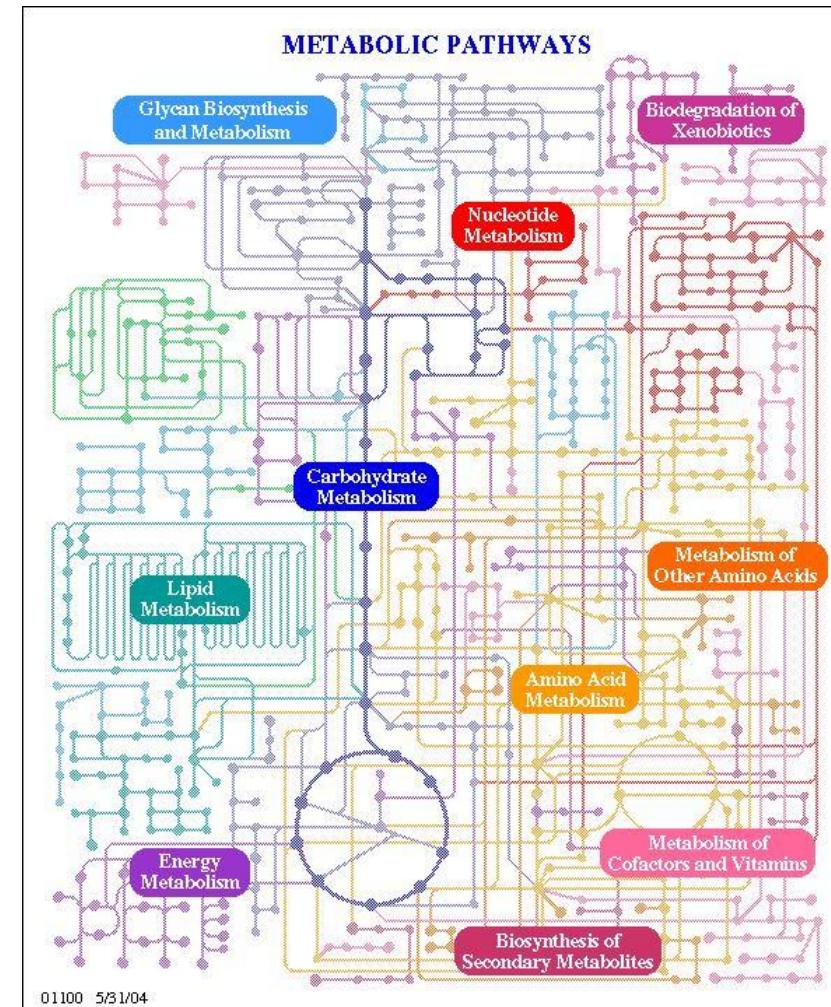
From the new paradigm, biological systems are complex networks of myriad of pathways, many of them interconnected



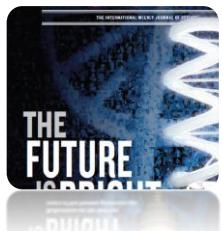


Biosynthesis  
pathways,

# Integrative Biology

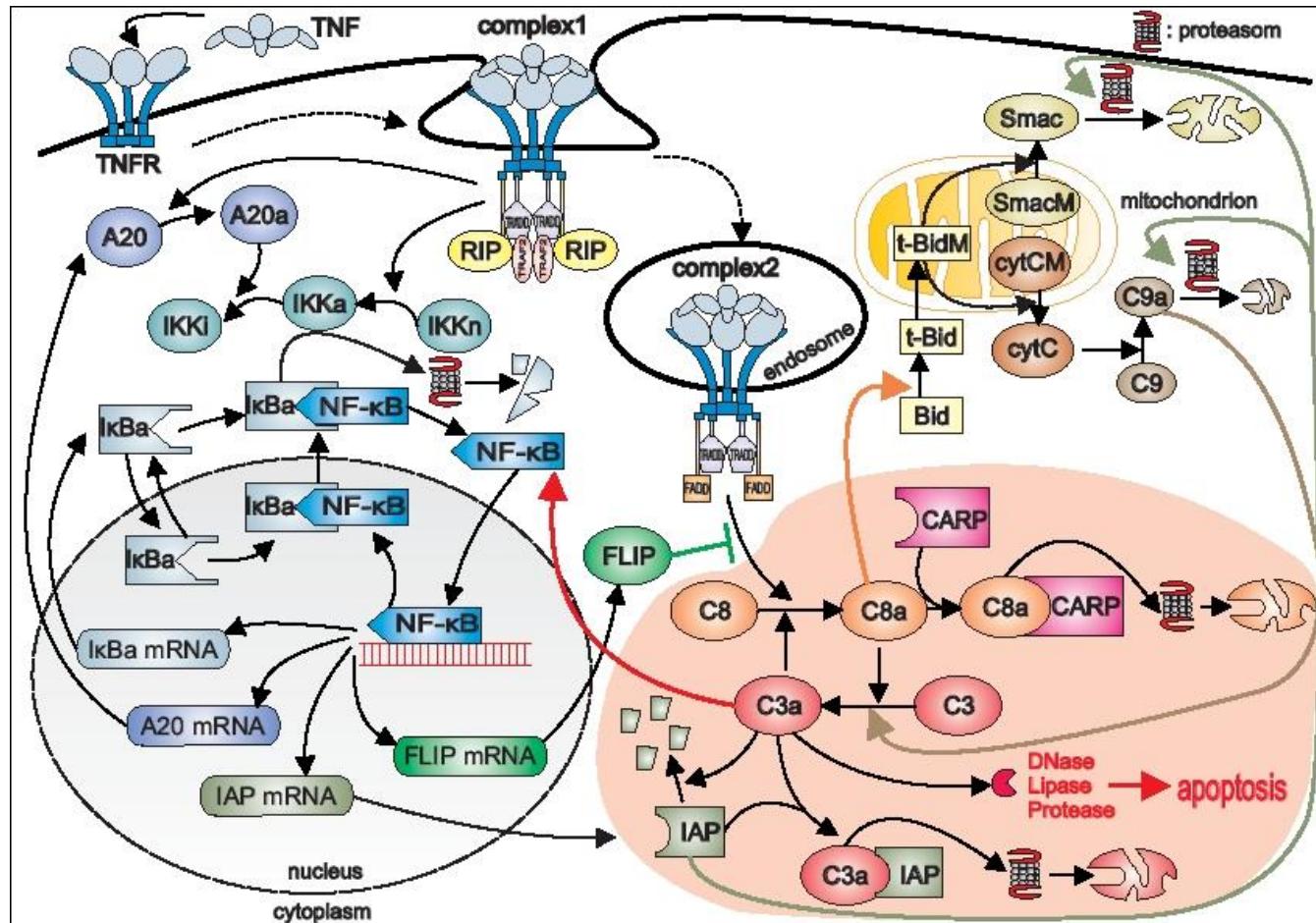


Metabolic pathway



# Integrative Biology

Signal transduction pathways,

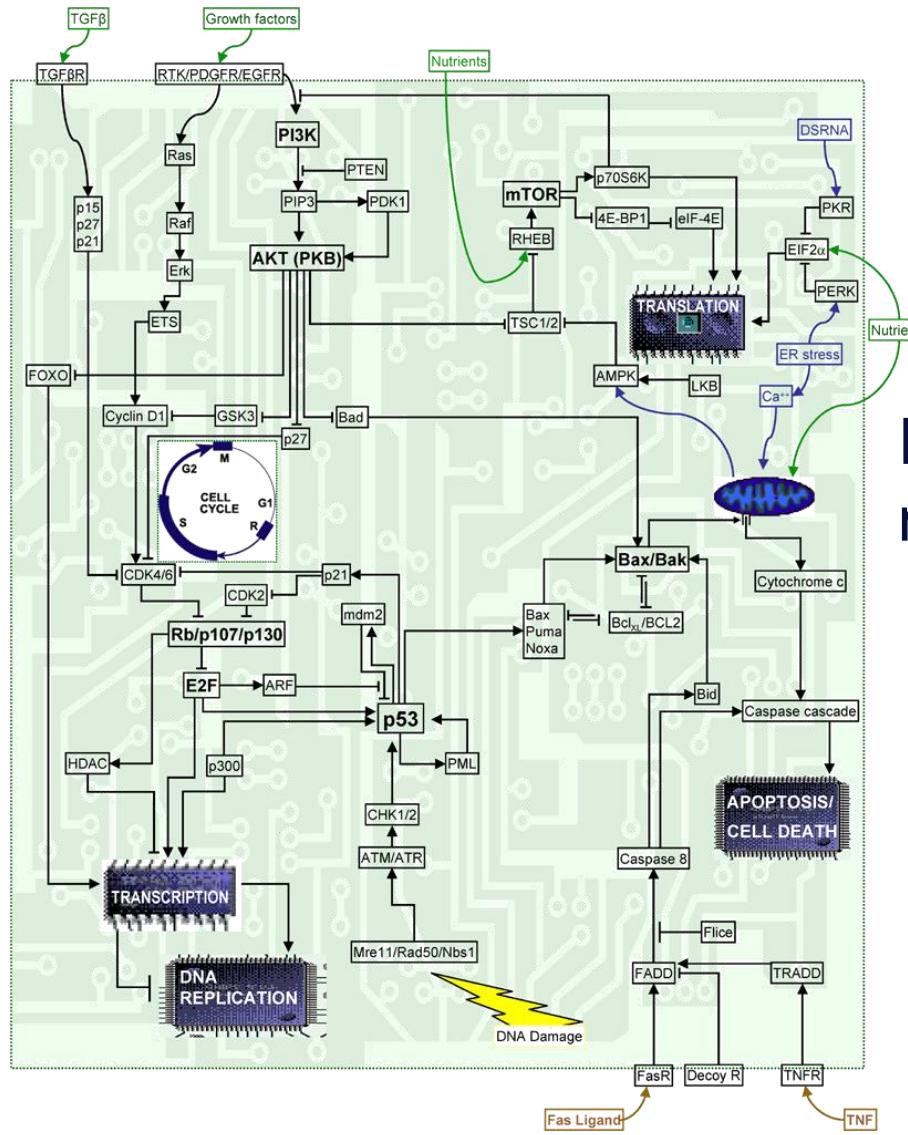


Signal transduction pathways



# Integrative Biology

Regulatory networks.

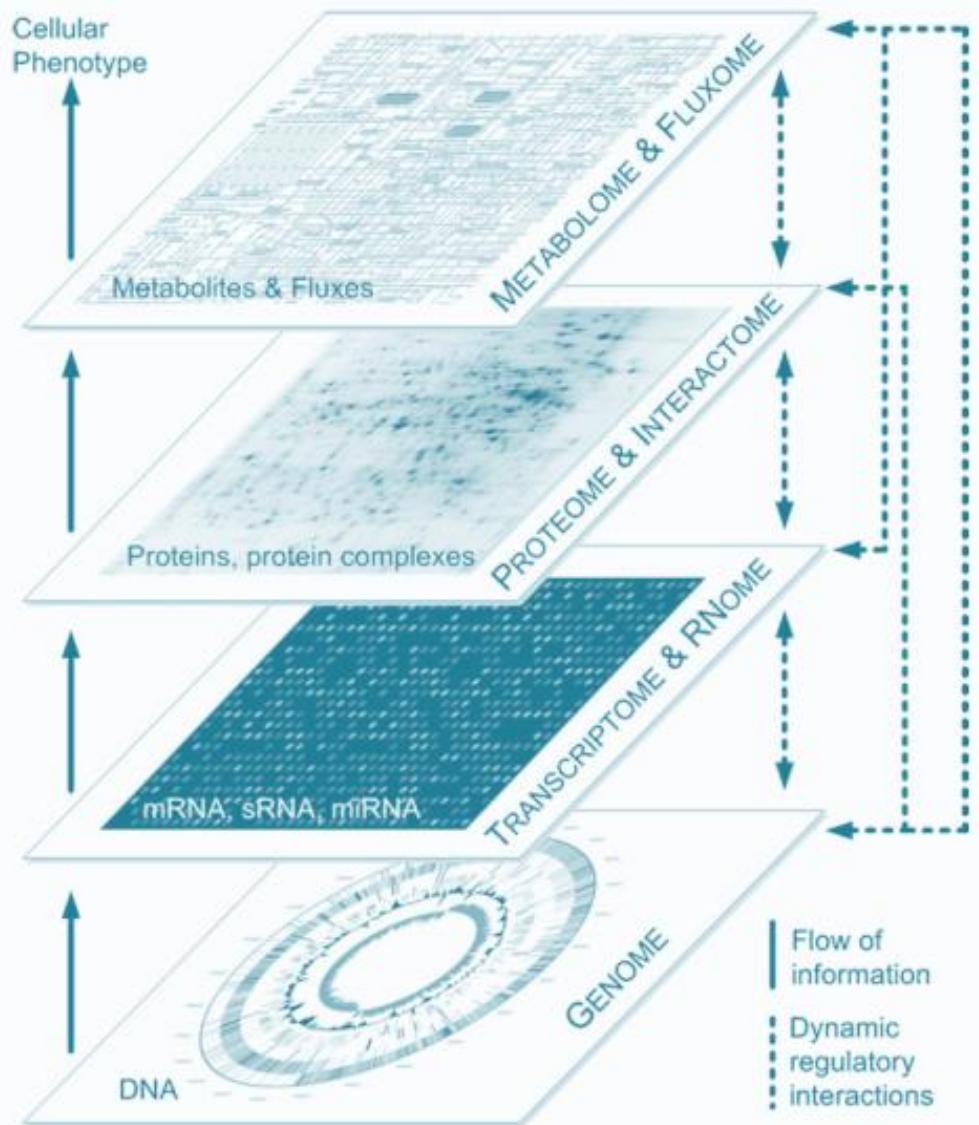


Regulatory  
networks



# Integrative Biology

Kohlsted et al. 2010



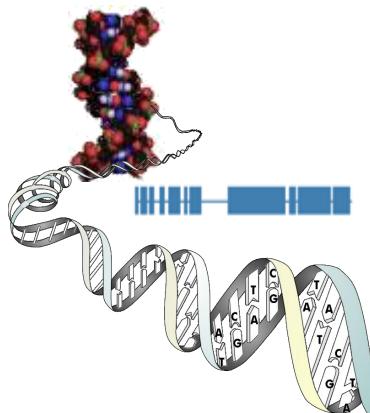


## The HGP has had a profound consequence in the **conceptualization of biological systems**

Era of the organism reconstruction: synthetic approach, interdisciplinary and big research teams to explain the emergent properties of biological systems

### 20th century biology

Reductionist approach  
(Experiments)



Biological System  
(Organism)



Construction blocks  
(Genes/Molecules)

Synthetic and interdisciplinary approach

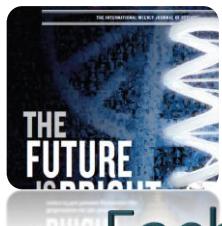
(Bioinformaticians, Biologists, Statisticians, Mathematicians, Biochemists, Physicist, Medical Doctors,...)



### 21st century biology



# Genome sequencing



# Why sequencing a genome?

Each genome sequence is a **treasure trove** containing an endless and invaluable source of biological information

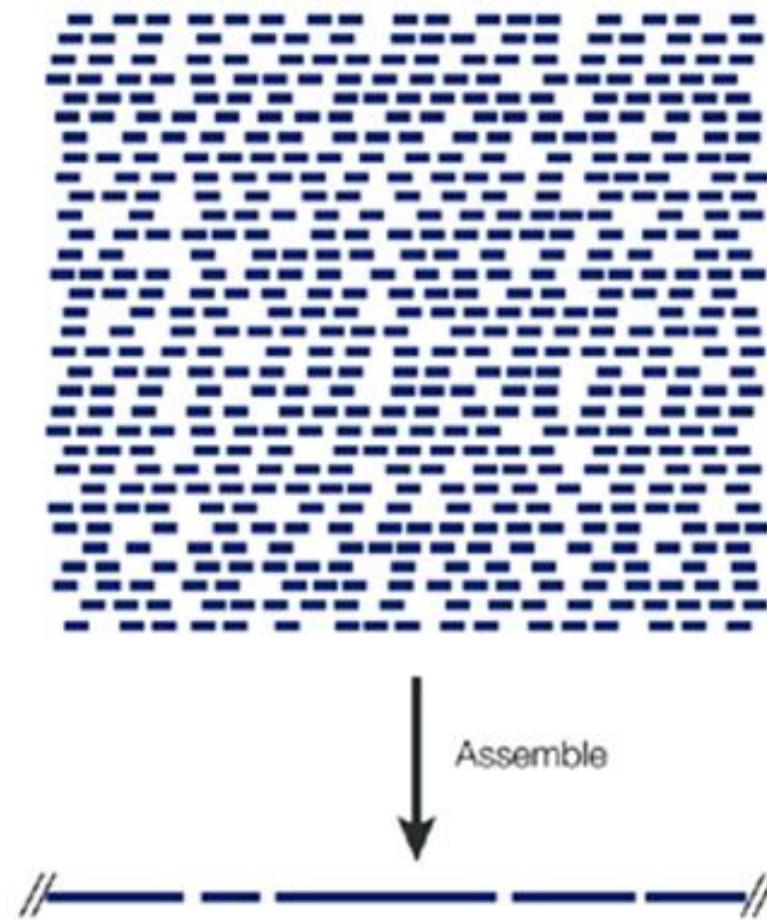
- Knowledge of the number, structure of genes and regulatory (functional) regions
- Basic Principles on the organization of the organism (functional classes, ...)
- Learn basic functions of genes conserved in different species (molecular biology lexicon)
- Chromosomal organization
- Genome evolution (conservation of gene order, sequence evolution)
- Genome variation (Population genomics)
- Association studies
- Expression analysis
- Integrative genomics (System biology)
- Applied genomics (Personalized medicine, Pharmacogenomics, Nutrigenomics, Agrigenomics, Conservation biology, Bioremediation,...)
- New areas of enquiry (Metagenomics, gene regulation, life evolution, human diaspora, medical genomics, ethics, law, ...)

We look at the forest, not a particular tree



## The pervasive assembly puzzle

How to map millions or billions of short read fragments onto a genome?





# The assembly puzzle

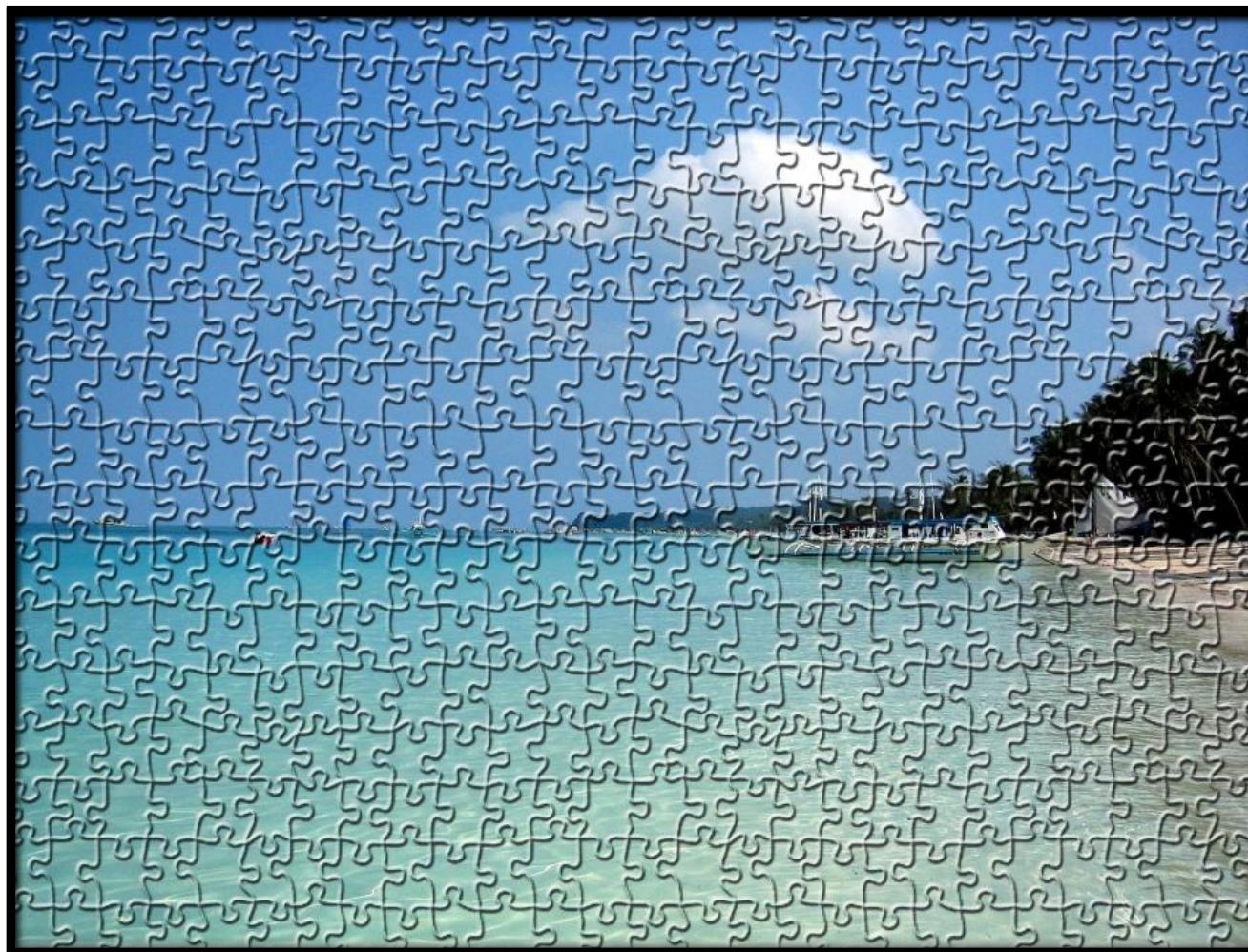
## Complex pattern (prokaryote genomes)





## The assembly puzzle

Simple pattern (genomes with high amount of repetitive sequence)





# Genome sequence metrics

Redundancy = Fold coverage

$$FC = \frac{N \cdot L}{G}$$

N = number of reads

L = mean read length

G = genome size

10 x is considered high quality

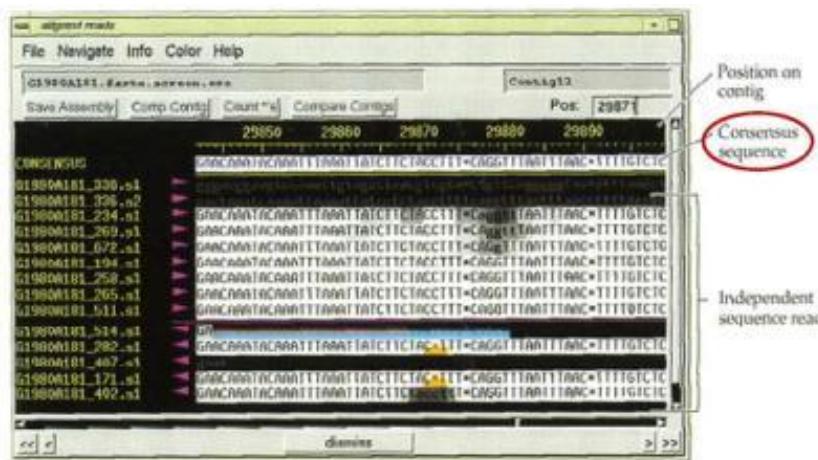
Base quality => phred score (Q)

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

P = Probability of calling a wrong base

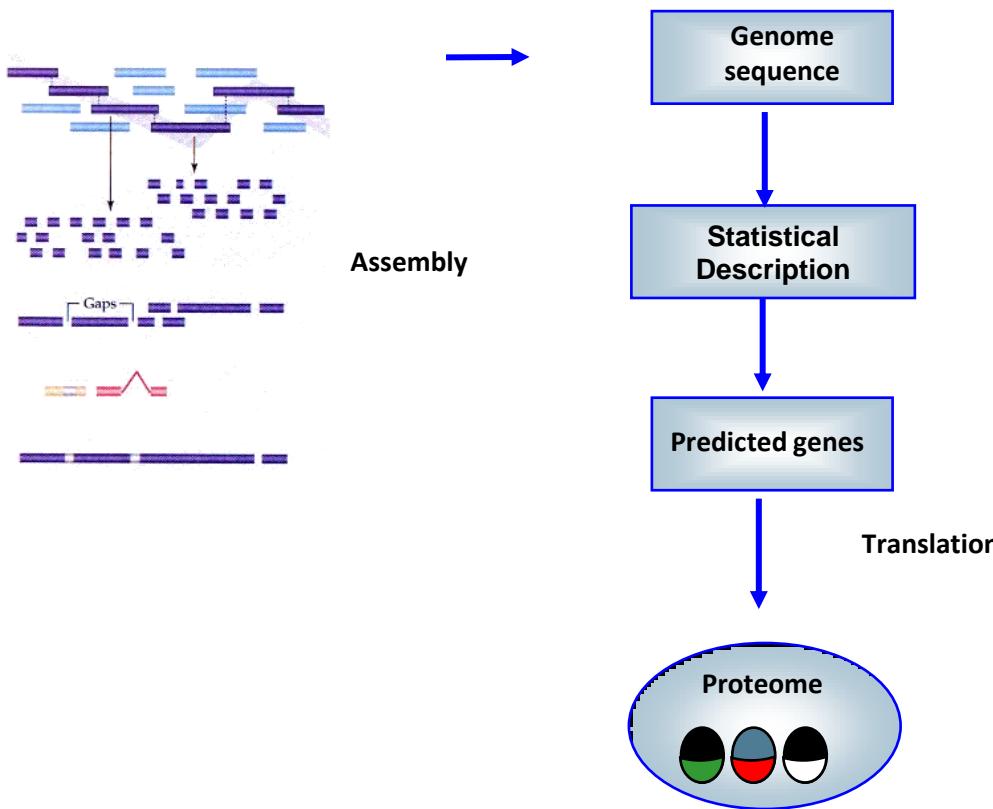
Q = 20 draft sequence (P = 0.01)

Q = 40 finished sequence (P = 0.0001)





# Steps of genome analysis



- DNA annotation and functional genomics
- Expression
- Proteins
- Molecular Evolution
- Genome variations and genotype-phenotype association studies
- System Biology



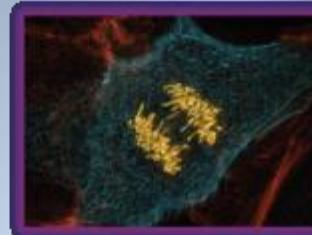
# The technological explosion



## Technology Advances Drive Science



Astronomy



Cell Biology

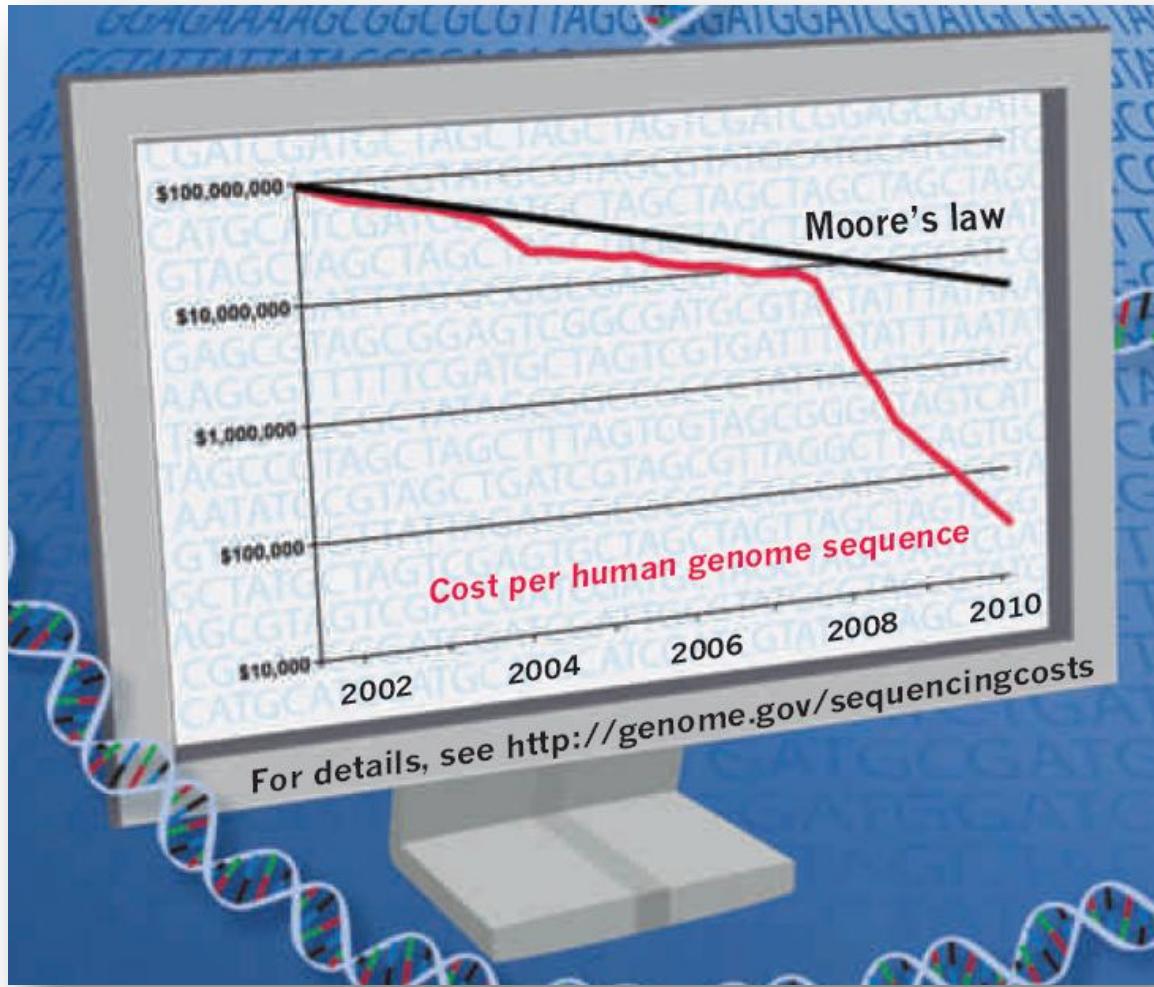


Radiology



Genomics

# The technological explosion

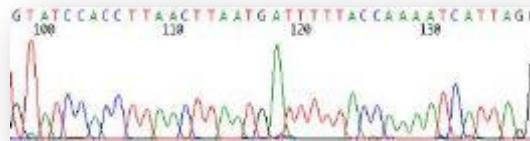




# Sequencing technologies



- Sanger (the *Gold Standard*) ([http://www.wiley-vch.de/books/sample/3527320903\\_c01.pdf](http://www.wiley-vch.de/books/sample/3527320903_c01.pdf)) Capillary electrophoresis. Reads (trace files) of 500-700 bases.



- Next Generation Technologies (starting in 2005):

## Massively parallel sequencing (\$1000 genome' program NHGRI)

- Dramatic increase of sequence output
- Significant decrease in the length of reads
- Decrease in the accuracy of calling bases -> Very different error profiles depending on the technological platform

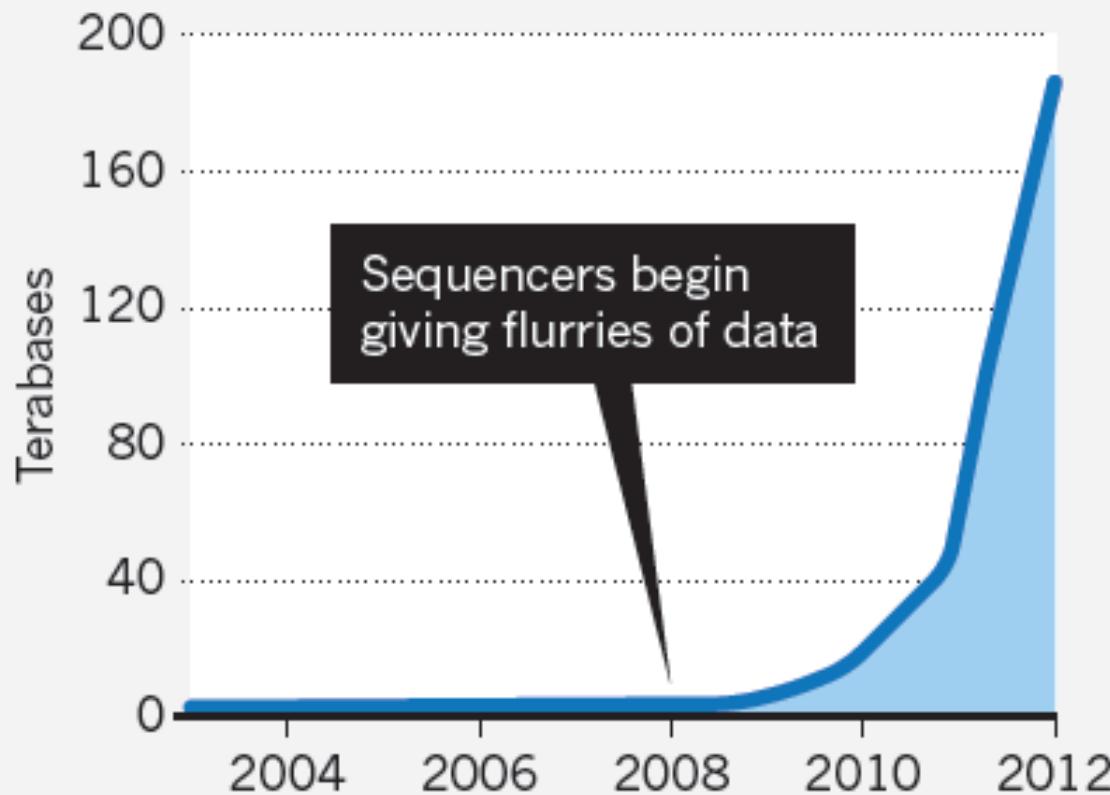


Analytical difficulties (profound changes in data analysis pipelines) → Revitalization of bioinformatics



## DATA EXPLOSION

The amount of genetic sequencing data stored at the European Bioinformatics Institute takes less than a year to double in size.





20 petabytes ~ 6 millones HG

# The European Bioinformatics Institute

Part of the European Molecular Biology Laboratory

EMBL-EBI provides freely available data from life science experiments, performs basic research in computational biology and offers an extensive user training programme, supporting researchers in academia and industry.

## Explore the EBI:

 Examples: blast, keratin, bf1 

## Press release



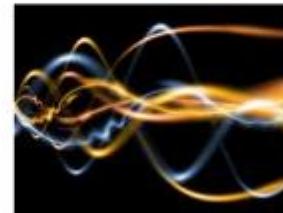
Bioinformatics embraces Semantic Web technologies

EMBL-EBI's new Resource Description Framework (RDF) platform provides access to bioinformatics resources that support Semantic Web technologies.



Functional genetic variation in humans: comprehensive map published

GEUVADIS project presents the largest-ever dataset linking human genomes to gene activity at the level of RNA.



It's not just noise  
EMBL scientists discern key gene expression patterns in the human genome

## Popular

- |          |          |
|----------|----------|
| Services | Jobs     |
| Research | Visit us |
| Training | EMBL     |
| News     | Contacts |

## Events

- [From Genome Sequencing to Gene Function - Izmir, Turkey](#)  
Oct 26 2013  
Registration deadline: Oct 25 2013
- [Managing and Exploring Next Generation Sequencing Data](#)  
Oct 29 2013 -Oct 30 2013  
Registration deadline: Oct 25 2013
- [Understanding 'omics data \(joint dixa - Genedata workshop\) - Basel, Switzerland](#)  
Nov 19 2013 -Nov 21 2013  
Registration deadline: Oct 27 2013
- [Mosquito Informatics \(INFRAVEC\)](#)  
Feb 5 2014 -Feb 6 2014  
Registration deadline: Nov 30 2013
- [Agricultural-Omics](#)  
Feb 17 2014 -Feb 21 2014  
Registration deadline: Dec 20 2013
- [Next Generation Sequencing Workshop](#)  
Mar 3 2014 -Mar 6 2014  
Registration deadline: Jan 3 2014
- [EMBO Practical Course on Metabolomics Bioinformatics for Life Scientists](#)  
Mar 17 2014 -Mar 21 2014  
Registration deadline: Jan 17 2014
- [Micro B3 Marine Metagenomics Bioinformatics](#)

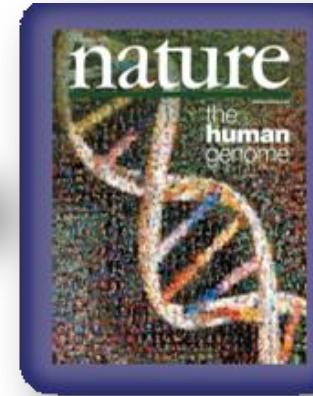
## Research infrastructures

EMBL-EBI is a pivotal partner in [ELIXIR](#), the European life sciences infrastructure for biological information... as part of the European Strategy on Research Infrastructures (ESRF1)... On



## The technological explosion

~\$1,000,000,000



Cost: 1 million fold lower!!!

←————— Today

~\$1,000

*“The \$1000 Genome”*



# Genome Sequencing as a “Commodity”

**Sherlock Holmes was an amateur.**



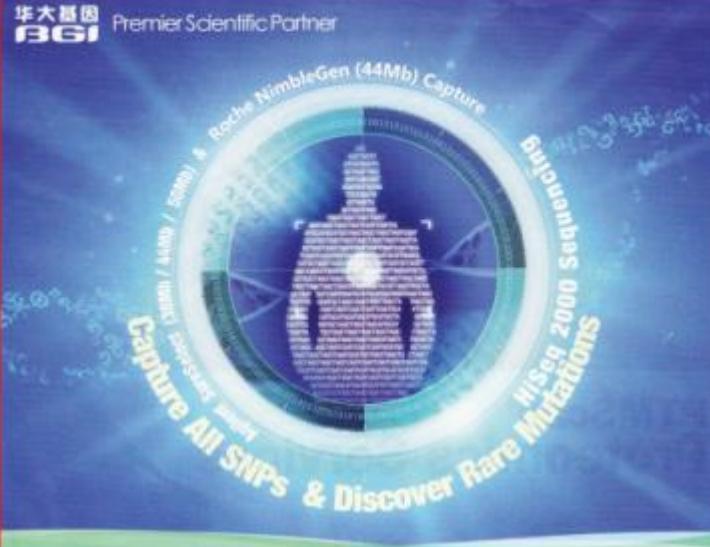
**SPECIAL PRICING \$4,998** Human Whole Genome Sequencing & Functional Interpretation (min. 10 genomes)

Investigating a genetic disease? We're the genome detectives to call. As experts in the functional interpretation of human genomes, we've built a state-of-the-art pipeline to rapidly annotate and thoroughly compare up to 300 whole genomes or exomes at once—to quickly track down the variants, genes, and pathways that govern disease. Starting with tissue samples, we deliver analyzed data, a shortlist of suspects, and powerful software to let you close the case in record time.

We can help you identify the variants, genes, and pathways that characterize a genetic disease. Visit [www.knome.com/disease](http://www.knome.com/disease) or call (857) 453-3895 to learn more.

**Knome**  
From DNA to Discovery

**BGI** Premier Scientific Partner



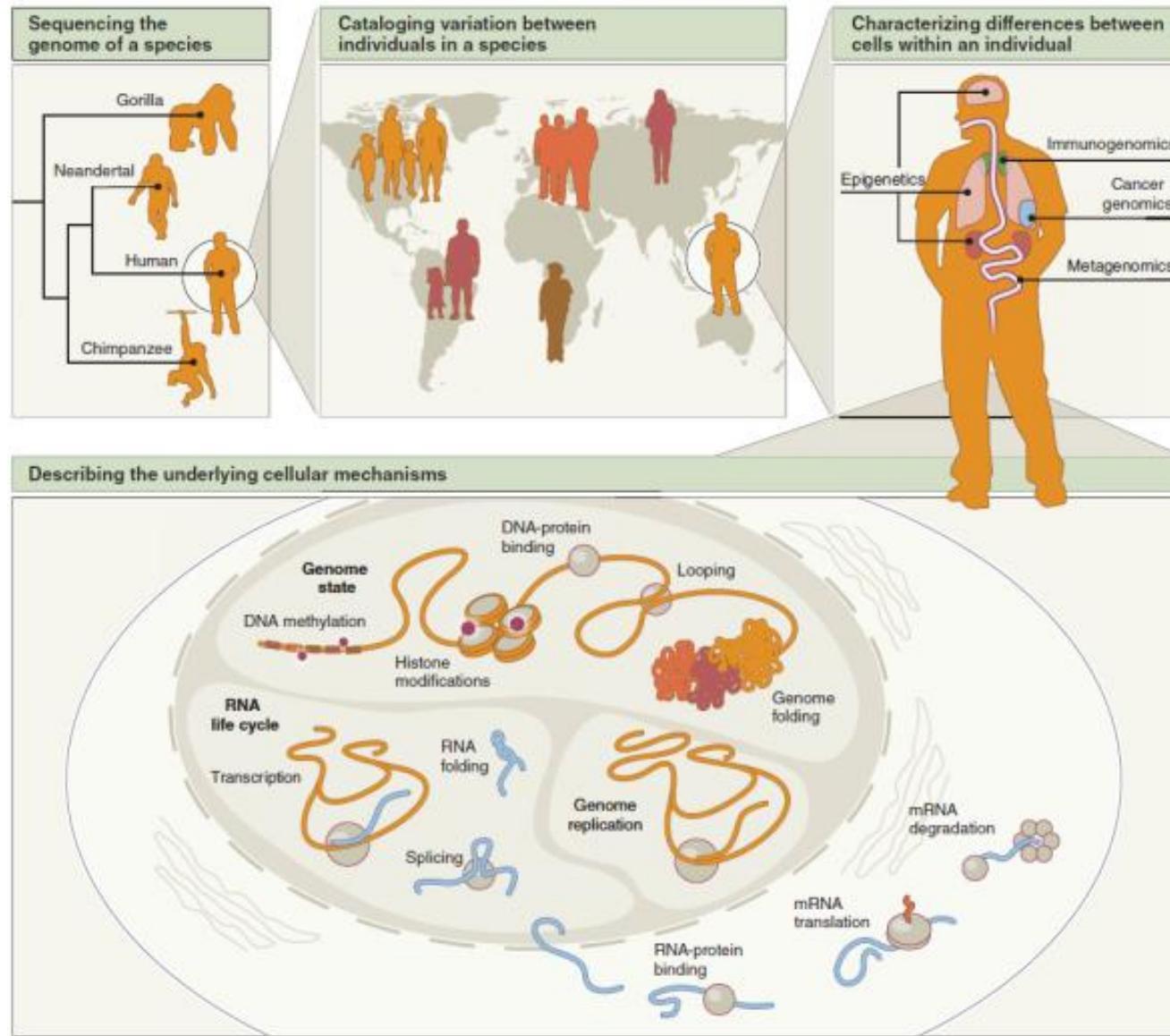
**Human Exome Sequencing Starting at \$999**

**Benefits**

- Target the most functionally relevant DNA sequences
- Capture both common and rare variants missed in traditional GWAS studies
- 150 next-generation sequencers assure rapid turnaround
- 1000 bioinformaticians generate high-quality, reliable data

America: (617) 900-2741 | [info@gnomelabs.com](mailto:info@gnomelabs.com) | Europe: +44 5000669 | [bsgi@bgi.com](mailto:bsgi@bgi.com) | [www.bgisequence.com](http://www.bgisequence.com)

# Road map of sequencing science





## Completed and ongoing genome projects

Complete Genome Projects: 7396

Archaeal: 234

Bacterial: 6851

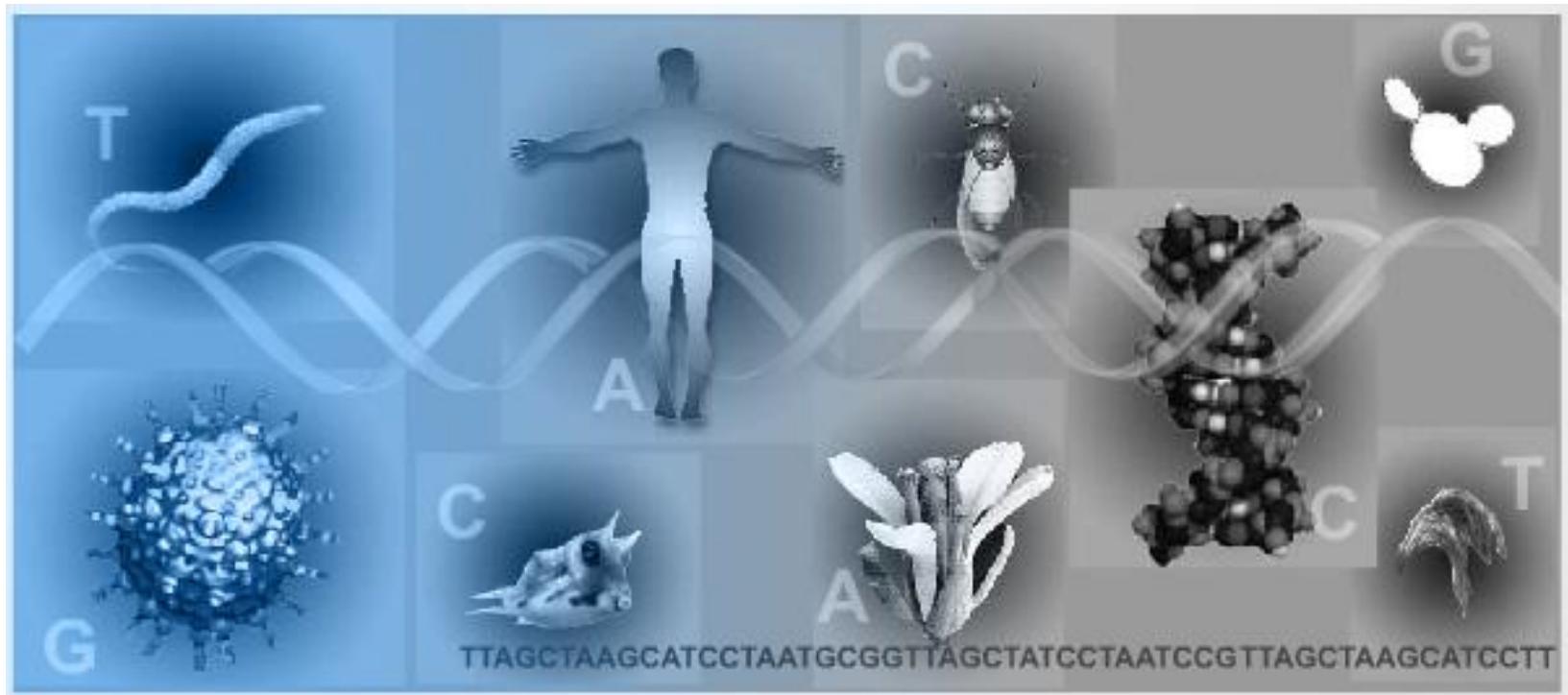
Eukaryal: 311



Finished: [2649](#)

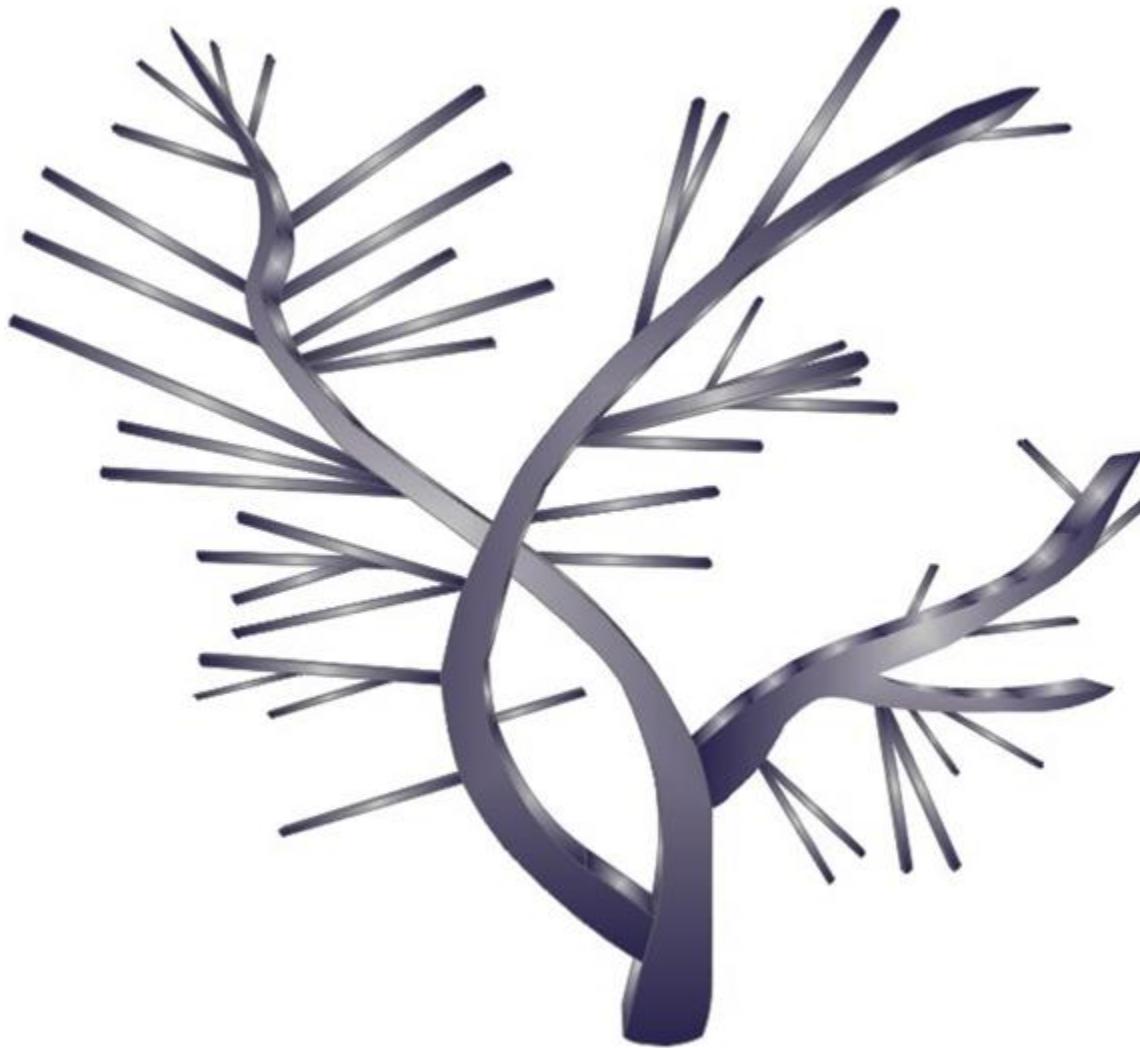


Permanent Draft: [4747](#)





## Genomas de especies





# Association Studies

## Variación fenotípica individual (incluye la predisposición a enfermedades)



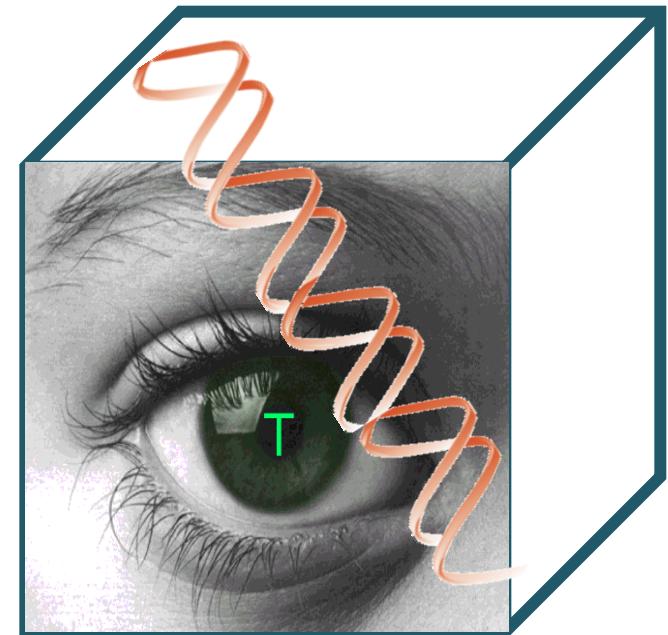
GENOMICS

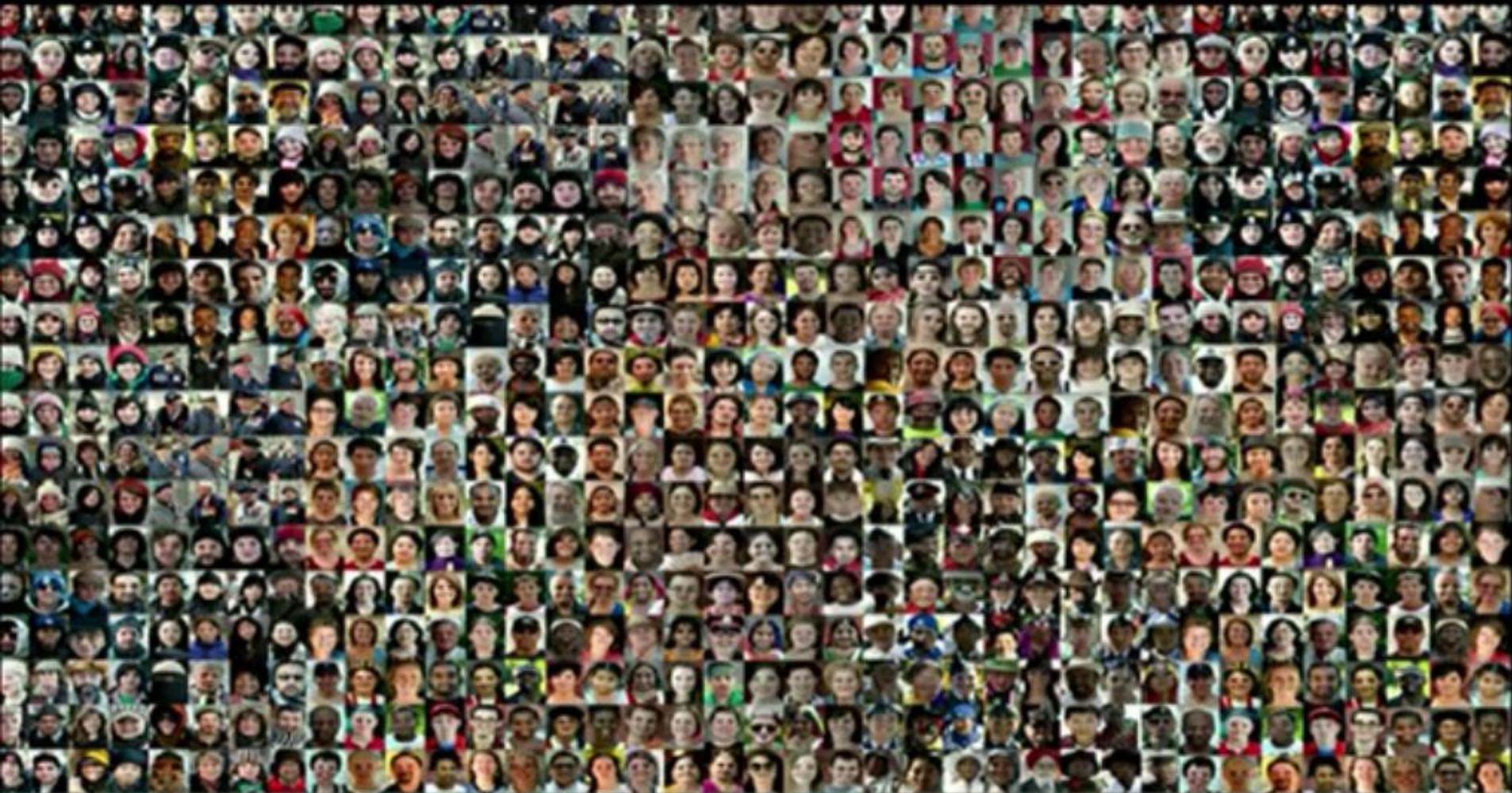
Fenotipo = Genotipo + Ambiente

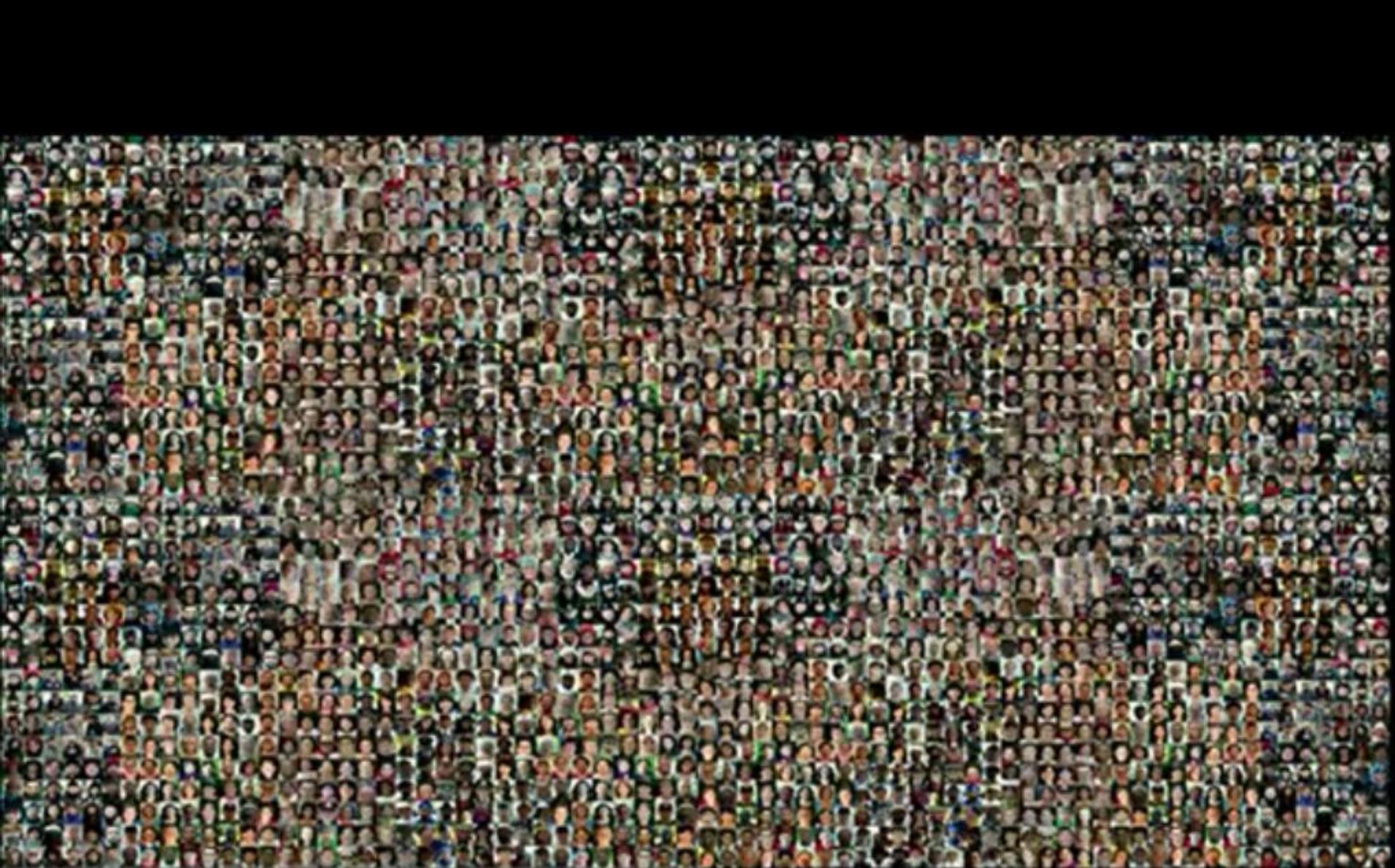


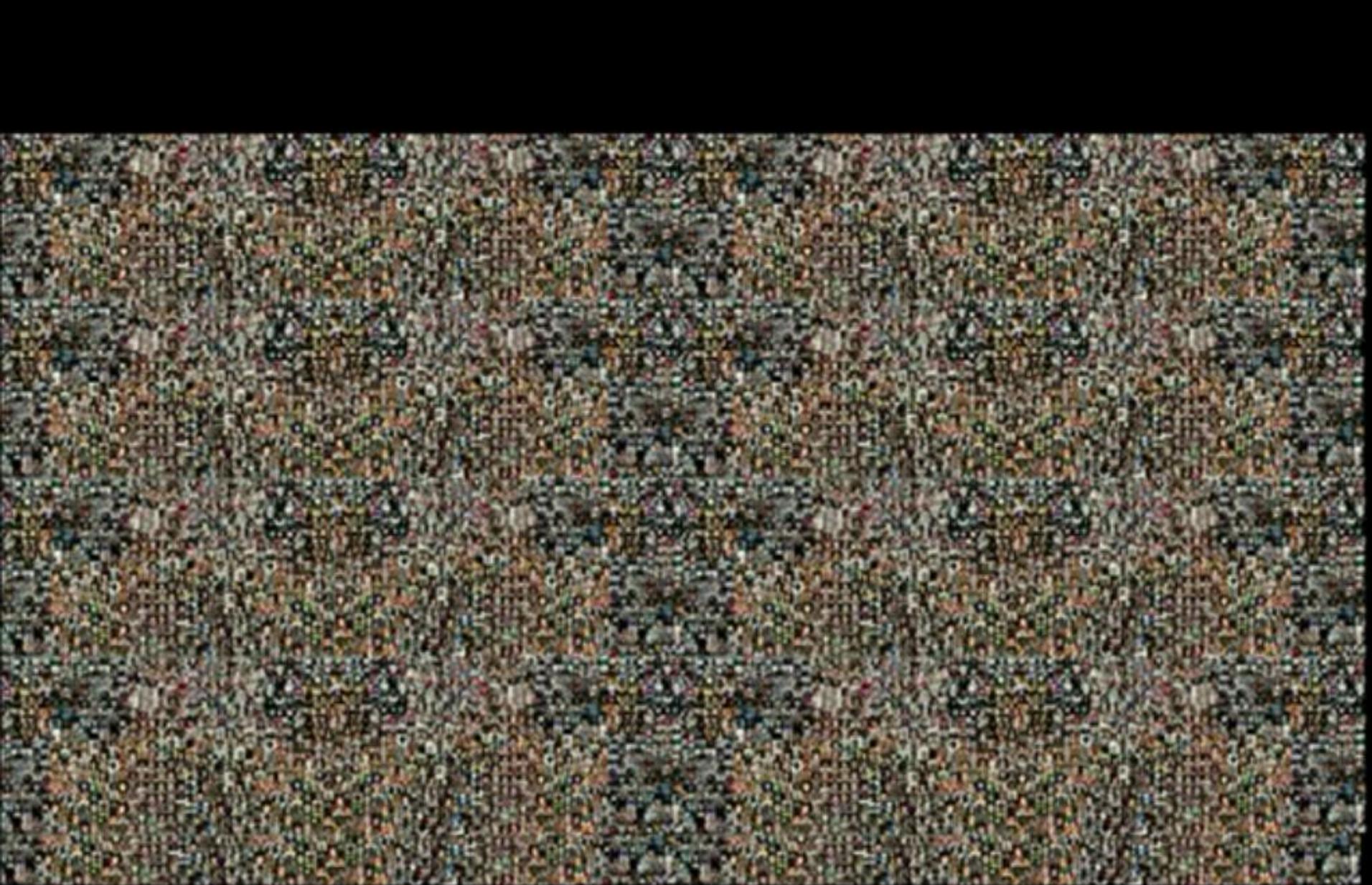
# The empirical space of human genetic variation:

## Data desideratum











# Association studies: Phenotypic effect of SNPs

Human genetic & phenotypic diversity database

	Genotype				Phenotype	
	SNP <sub>1</sub>	SNP <sub>2</sub>	SNP <sub>3</sub>	...	Disease 1	Trait i
Secuence individual 1	A/A	G/C	G/T		Healthy	$x_1$
Secuence individual 2	A/C	C/C	T/T		Cervical Cancer	$x_2$
...					...	...

Estimation phenotypic effect

G/C | G/T

Healthy

$x_1$

Cervical  
Cancer

$x_2$



## BioBanks: Studies of cohorts at a great scale

# The UK Biobank

A study of genes, environment and health

USA



- deCODE (Islandia)
- Estonia
- Germany
- Canada
- Japan
- China

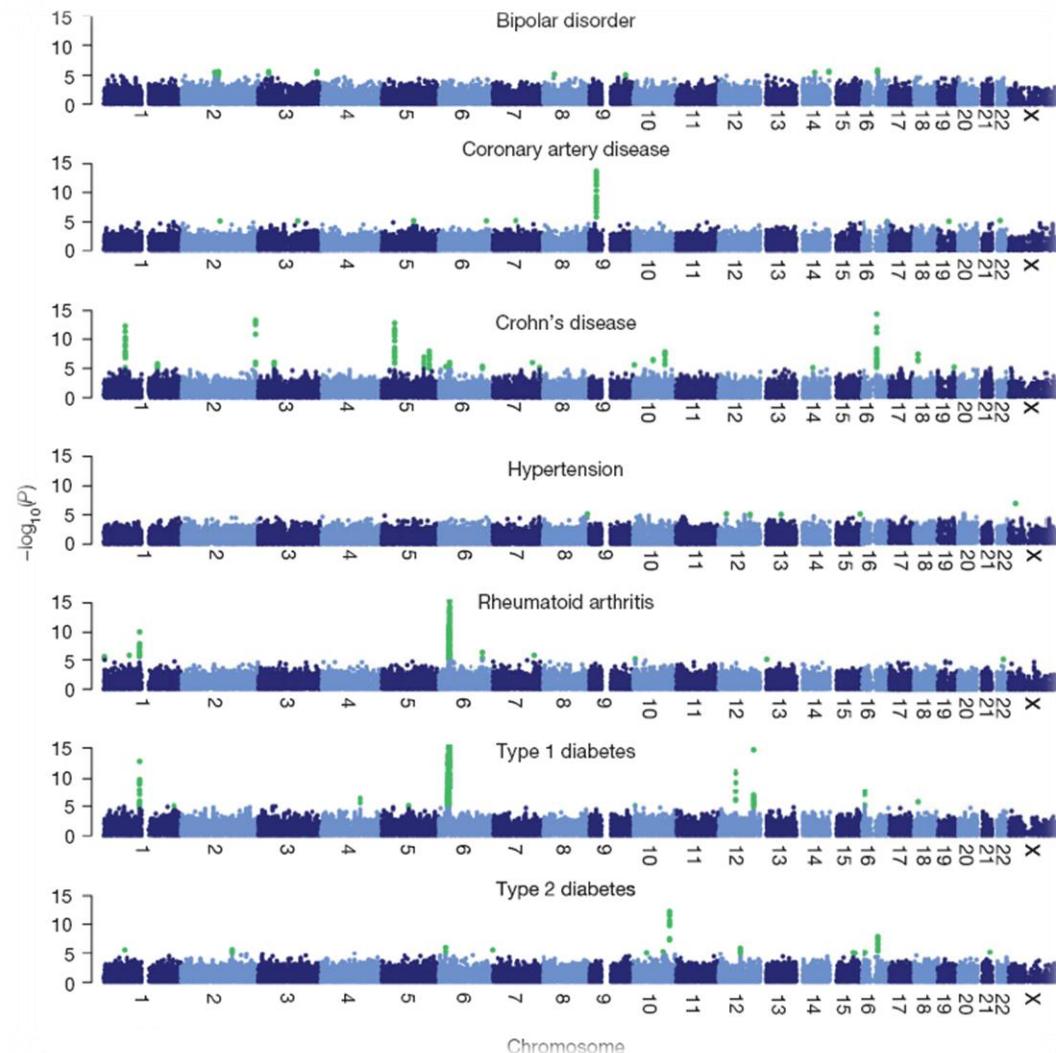


# Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls

The Wellcome Trust Case Control Consortium\*

## GWA Big Data problems

- Data bias
- False positives
- Data structure
- Security issues

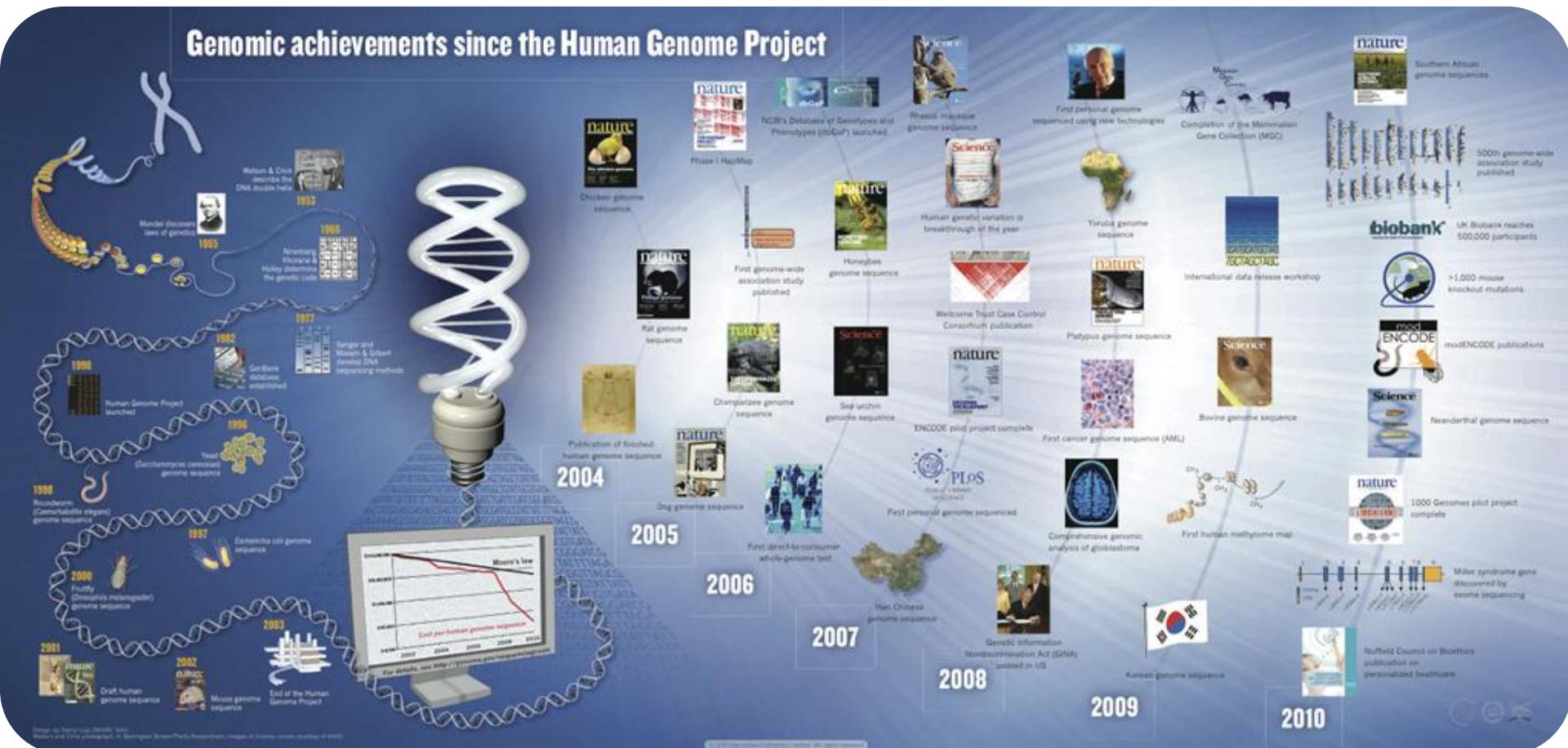




# The triumphal march of genomics

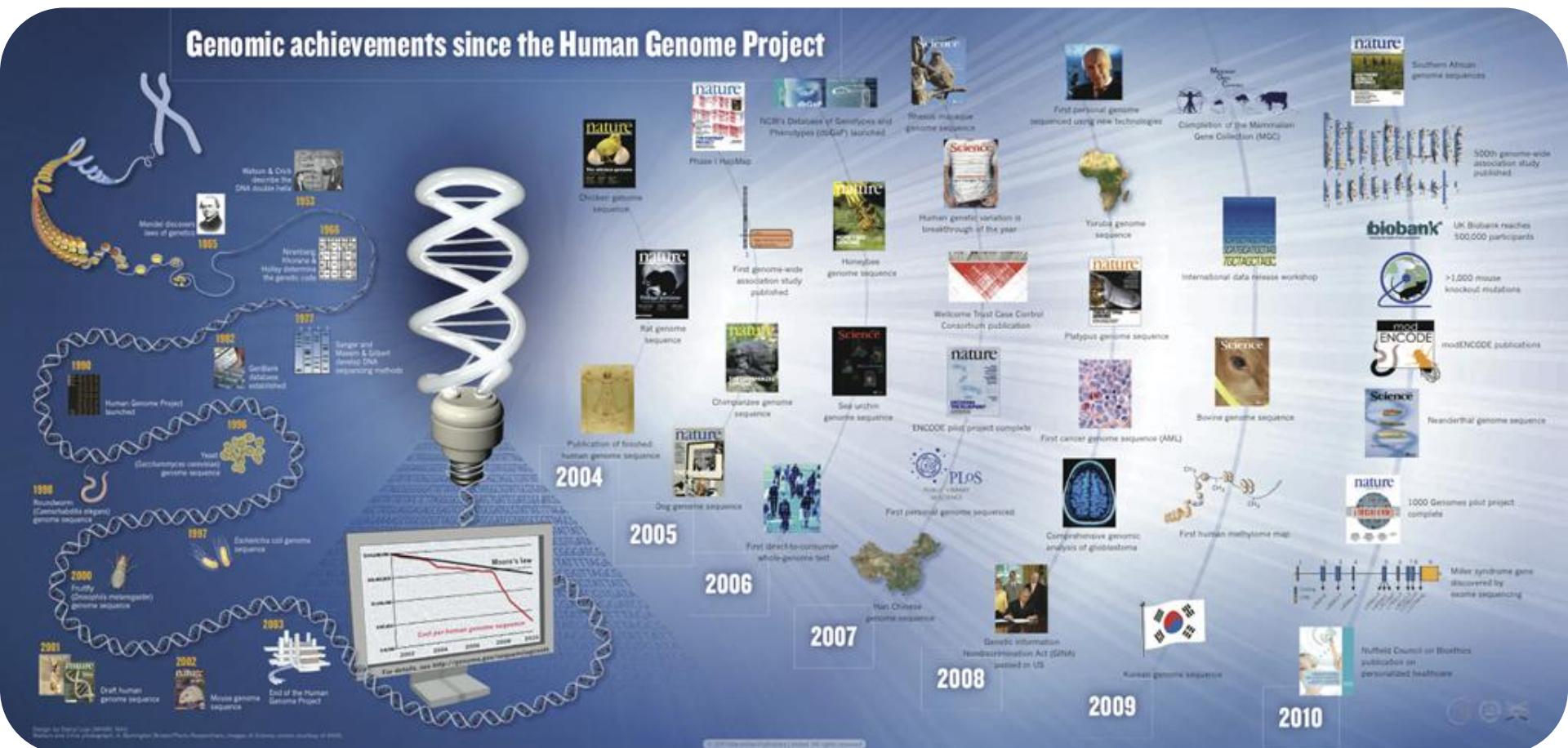


# The triumphal march of genomics: from the human genome to the 1000 human genomes



ED Green et al. *Nature* **470**, 204-213 (2011) doi:10.1038/nature09764

# Genomic achievements since the Human Genome Project



ED Green *et al.* *Nature* **470**, 204-213 (2011) doi:10.1038/nature09764

nature



# Genome science challenges

**BOX 3**

## Bioinformatics and computational biology



The major bottleneck in genome sequencing is no longer data generation—the computational challenges around data analysis, display and integration are now rate limiting. New approaches and methods are required to meet these challenges.

**Data analysis.** Computational tools are quickly becoming inadequate for analysing the amount of genomic data

that can now be generated, and this mismatch will worsen. Innovative approaches to analysis, involving close coupling with data production, are essential.

**Data integration.** Genomics projects increasingly produce disparate data types (for example, molecular, phenotypic, environmental and clinical), so computational approaches must not only keep pace with the volume of genomic data, but also their complexity. New integrative methods for analysis and for building predictive models are needed.

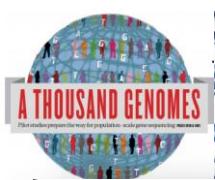
**Visualization.** In the past, visualizing genomic data involved indexing to the one-dimensional representation of a genome. New visualization tools will need to accommodate the multidimensional data from studies of molecular phenotypes in different cells and tissues, physiological states and developmental time. Such tools must also incorporate non-molecular data, such as phenotypes and environmental exposures. The new tools will need to accommodate the scale of the data to deliver information rapidly and efficiently.

**Computational tools and infrastructure.** Generally applicable tools are needed in the form of robust, well-engineered software that meets the distinct needs of genomic and non-genomic scientists. Adequate computational infrastructure is also needed, including sufficient storage and processing capacity to accommodate and analyse large, complex data sets (including metadata) deposited in stable and accessible repositories, and to provide consolidated views of many data types, all within a framework that addresses privacy concerns. Ideally, multiple solutions should be developed<sup>105</sup>.

**Training.** Meeting the computational challenges for genomics requires scientists with expertise in biology as well as in informatics, computer science, mathematics, statistics and/or engineering. A new generation of investigators who are proficient in two or more of these fields must be trained and supported.

## The major bottleneck in genome sequencing is computational challenges around data analysis, display and integration

- **Data analysis:** Keep pace with the volume of genomic data
- **Data integration:** Keep pace with the complexity of genomic data.
- **Visualization:** New visualization tools will need to accommodate the multidimensional data from studies of molecular phenotypes in different cells and tissues, physiological states and developmental time.
- **Computational tools and infrastructure:**
  - Robust, well-engineered software that meets the distinct needs of genomic and non-genomic scientists.
  - Adequate computational infrastructure: sufficient storage and processing capacity to accommodate and analyse large, complex data sets deposited in stable and accessible repositories.
- **Training:** A new generation of investigators proficient in two or more of fields of biology, informatics, computer science, mathematics, statistics and/or engineering must be trained and supported.



## The 1000 Genomes Project

### ARTICLE

#### A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium\*

<http://www.nature.com/nature/journal/v467/n7319/pdf/nature09534.pdf>

The 1000 Genomes Project aims to provide a deep characterization of human genome sequence variation as a foundation for investigating the relationship between genotype and phenotype. Here we present results of the pilot phase of the project, designed to develop and compare different strategies for genome-wide sequencing with high-throughput platforms. We undertook three projects: low-coverage whole-genome sequencing of 179 individuals from four populations; high-coverage sequencing of two mother-father-child trios; and exon-targeted sequencing of 697 individuals from seven populations. We describe the location, allele frequency and local haplotype structure of

### ARTICLE

doi:10.1038/nature11632

#### An integrated map of genetic variation from 1,092 human genomes

The 1000 Genomes Project Consortium\*

By characterizing the geographic and functional spectrum of human genetic variation, the 1000 Genomes Project aims to build a resource to help to understand the genetic contribution to disease. Here we describe the genomes of 1,092 individuals from 14 populations, constructed using a combination of low-coverage whole-genome and exome sequencing. By developing methods to integrate information across several algorithms and diverse data sources, we provide a validated haplotype map of 38 million single nucleotide polymorphisms, 1.4 million short insertions and deletions, and more than 14,000 larger deletions. We show that individuals from different populations carry different profiles of rare and common variants, and that low-frequency variants show substantial geographic differentiation, which is further increased by the action of purifying selection. We show that evolutionary conservation and coding consequence are key determinants of the strength of purifying selection, that rare-variant load varies substantially across biological pathways, and that each individual contains hundreds of rare non-coding variants at conserved sites, such as motif-disrupting changes in transcription-factor-binding sites. This resource, which captures up to 98% of accessible single nucleotide polymorphisms at a frequency of 1% in related populations, enables analysis of common and low-frequency variants in individuals from diverse, including admixed, populations.

1000 genome project: The 1000 Genomes Project Consortium. 2012. An integrated map of genetic variation from 1,092 human genomes. [Nature 491: 56-65](https://doi.org/10.1038/nature11632)

# 1000 genomes Big Data problems

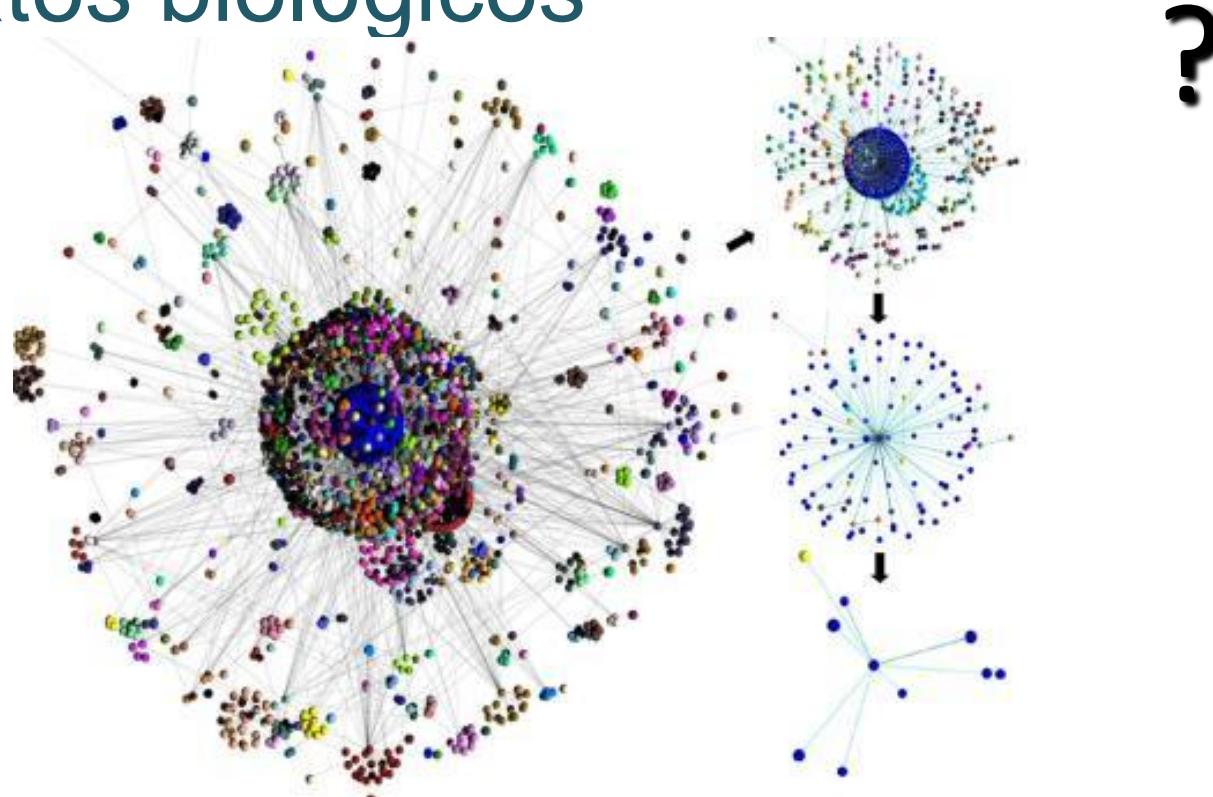
- Data storage (local vs cloud)
- Data computing
- Data delivery
- Security issues



# The foreseeable future of genome science



# ¿Qué gran panorámica emergirá de la montaña de datos biológicos



- a. La complejidad no es reducible
- b. Nuevos principios generales de organización de lo biológico



## NHGRI Plan: Charting a course for genomic medicine from base pairs to bedside

genome.gov  
National Human Genome Research Institute  
National Institutes of Health

Research Funding | Research at NHGRI | Health | Education | Issues in Genetics | Newsroom | Careers & Training | About | For You | [Facebook](#) [Twitter](#) [YouTube](#)

[Home](#) > [About](#) > [Long-Range Planning](#) > Charting a course for genomic medicine from base pairs to bedside: The Strategic Plan

**Long-Range Planning**

- Charting a course for genomic medicine from base pairs to bedside: The Strategic Plan**
- Event: A Decade with the Human Genome Sequence: Charting a Course for Genomic Medicine
- Past Long-Range Planning
- Topics: NHGRI 2008-2011 Planning Process
- White Papers: The 2008-2011 Planning Process
- Workshops: NHGRI 2008-2011 Planning Process

**The Strategic Plan**

**Charting a course for genomic medicine from base pairs to bedside**

**On February 10, 2011, Nature magazine published the National Human Genome Research Institute's (NHGRI) strategic plan for the future of human genome research called *Charting a course for genomic medicine from base pairs to bedside*. This strategic vision was developed in consultation with leading genome researchers over more than two years and is intended to inspire many to contribute to advancing genomic understanding, especially as other National Institutes of Health (NIH) institutes and centers focus genomic technologies on the diseases they study.**

To celebrate the 10th anniversary of the first analysis of the draft human genome, and the launch of the new strategic vision for the field of genomics, NHGRI sponsored a symposium with leading thinkers in the field of genome research — including all three directors in the history of the National Human Genome Research Institute — who gathered on the campus of the National Institutes of Health on February 11, 2011, to consider the future of their field.

**Follow the links below to the full strategic plan, symposium information, the planning process that developed the plan and related information.**

- **Strategic Plan:** [Charting a course for genomic medicine from base pairs to bedside](#)
- February 10, 2011
- **Symposium:** [A Decade with the Human Genome Sequence: Charting a Course for Genomic Medicine](#)
- **Strategic Plan Press Release:** [NHGRI charts course for the next phase of genomics research](#)
- [NHGRI Long-Range Planning](#)

To view the PDF document(s) on this page, you will need Adobe Reader.

[Top of page](#)

Last Updated: March 23, 2012

<http://www.genome.gov/Pages/About/Planning/2011NHGRIStrategicPlan.pdf>

**THE  
FUTURE  
IS BRIGHT**

Reflections on the first ten  
years of the human genomics age



Years of the human genomics age  
Reflections on the first ten

February 2011

NHGRI Published New Vision for Genomics

## PERSPECTIVE

doi:10.1038/nature09764

### Charting a course for genomic medicine from base pairs to bedside

Eric D. Green<sup>1</sup>, Mark S. Guyer<sup>2</sup> & National Human Genome Research Institute\*

There has been much progress in genomics in the ten years since a draft sequence of the human genome was published. Opportunities for understanding health and disease are now unprecedented, as advances in genomics are harnessed to obtain robust foundational knowledge about the structure and function of the human genome and about the genetic contributions to human health and disease. Here we articulate a 2011 vision for the future of genomics research and describe the path towards an era of genomic medicine.

Since the end of the Human Genome Project (HGP) in 2003 and the publication of a reference human genome sequence<sup>1,2</sup>, genomics has become a mainstay of biomedical research. The scientific community's foresight in launching this ambitious project<sup>3</sup> is evident in the broad range of scientific advances that the HGP has enabled, as shown in Fig. 1 (see rollfold). Optimism about the potential contributions of genomics for improving human health has been fuelled by new insights about cancer<sup>4–7</sup>, the molecular basis of inherited diseases (<http://www.ncbi.nlm.nih.gov/omim> and <http://www.genome.gov/GWASStudies>) and the role of structural variation in disease<sup>8</sup>, some of which have already led to new therapies<sup>9–11</sup>. Other advances have already changed medical practice (for example, microarrays are now used for clinical detection of genomic imbalances<sup>12</sup> and pharmacogenomic testing is routinely performed before administration of certain medications<sup>13</sup>). Together, these achievements (see accompanying paper<sup>14</sup>) document that genomics is contributing to a better understanding of human biology and to improving human health.

As it did eight years ago<sup>14</sup>, the National Human Genome Research Institute (NHGRI) has engaged the scientific community (<http://www.genome.gov/Planning>) to reflect on the key attributes of genomics (Box 1) and explore future directions and challenges for the field. These discussions have led to an updated vision that focuses on understanding human biology and the diagnosis, prevention and treatment of human disease, including consideration of the implications of those advances for society (but these discussions, intentionally, did not address the role of genomics in agriculture, energy and other areas). Like the HGP, achieving this vision is broader than what any single organization or country can achieve—realizing the full benefits of genomics will be a global effort.

quickly. Although genomics has already begun to improve diagnostics and treatments in a few circumstances, profound improvements in the effectiveness of healthcare cannot realistically be expected for many years (Fig. 2). Achieving such progress will depend not only on research, but also on new policies, practices and other developments. We have illustrated the kinds of achievements that can be anticipated with a few examples (Box 2) where a confluence of need and opportunities should lead to major accomplishments in genomic medicine in the coming decade. Similarly, we note three cross-cutting areas that are broadly relevant and fundamental across the entire spectrum of genomics and genomic medicine: bioinformatics and computational biology (Box 3), education and training (Box 4), and genomics and society (Box 5).

#### Understanding the biology of genomes

Substantial progress in understanding the structure of genomes has revealed much about the complexity of genome biology. Continued acquisition of basic knowledge about genome structure and function will be needed to illuminate further those complexities (Fig. 2). The contribution of genomics will include more comprehensive sets (catalogues) of data and new research tools, which will enhance the capabilities of all researchers to reveal fundamental principles of biology.

#### Comprehensive catalogues of genomic data

Comprehensive genomic catalogues have been uniquely valuable and widely used. There is a compelling need to improve existing catalogues and to generate new ones, such as complete collections of genetic variation, functional genomic elements, RNAs, proteins, and other biological

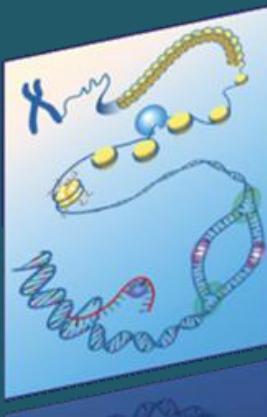


## Five domains of genomics research

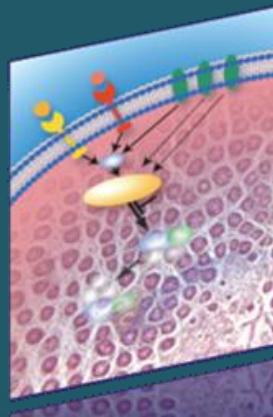
Understanding  
the Structure of  
Genomes



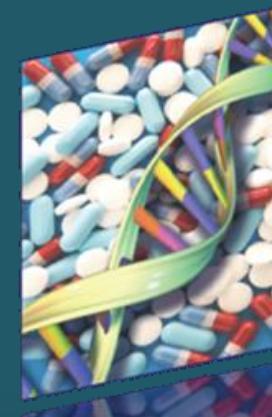
Understanding  
the Biology of  
Genomes



Understanding  
the Biology of  
Disease



Advancing  
the Science of  
Medicine



Improving the  
Effectiveness  
of Healthcare



**Basic  
Science**

**Translational  
Science**

**Implementation  
Science**



## Schematic representation of accomplishments across five domains of genomics research

# The challenges of big data in Genomics

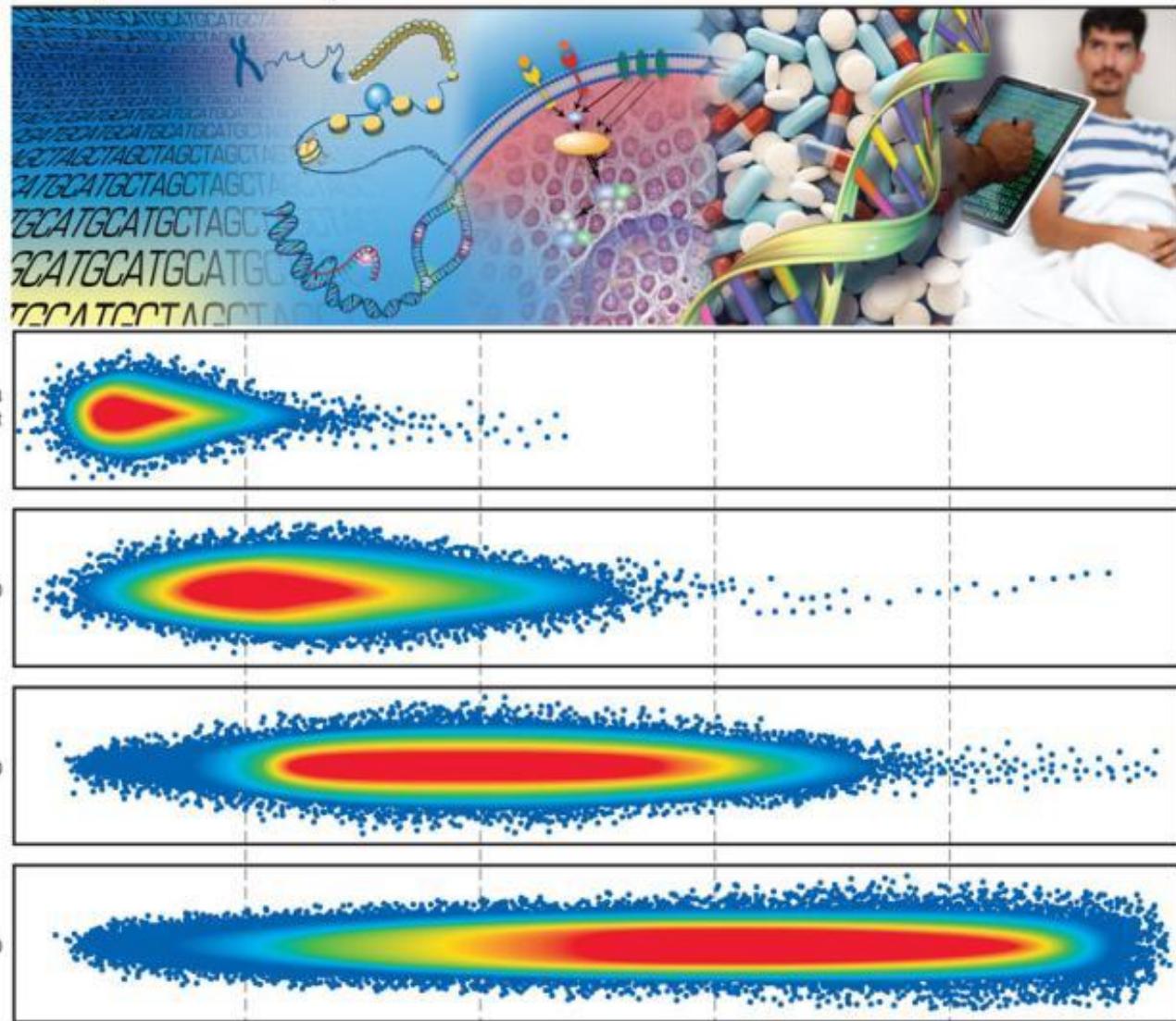
Understanding  
the structure of  
genomes

Understanding  
the biology of  
genomes

Understanding  
the biology of  
disease

Advancing  
the science of  
medicine

Improving the  
effectiveness of  
healthcare



E D. Green *et al.* *Nature* **470**, 204-213 (2011) doi:10.1038/nature09764

**nature**

*This is the greatest  
intellectual moment  
in history*



genome.gov



THE BRIGHT FUTURE  
OF HUMAN  
GENOMICS



# Readings & Videos

## Readings

- Perspective: [\*\*Charting a course for genomic medicine from base pairs to bedside.\*\*](#) 2011. Eric D. Green, Mark S. Guyer & National Human Genome Research Institute. Nature 470, 204–213 (10 February 2011).
- Review: [\*\*Initial impact of the sequencing of the human genome\*\*](#). Eric S. Lander. Nature 470, 187-197 (10 February 2011)
- Report economic study: [\*\*Economic Impact of the Human Genome Project\*\*](#). 2011. May 2011. Battelle technology partnership practice.
- [\*\*The \\$1,000 genome, the \\$100,000 analysis?\*\*](#). 2010. Elaine R Mardis. Genome Medicine 2010, 2:84

---

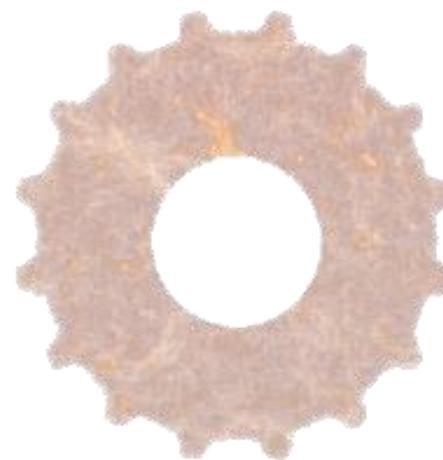
## Videos

- [\*\*Nobel Week Dialogue 2012. The Genetic Revolution and its Impact on Society\*\*](#)
- [\*\*The Genomic Landscape circa 2012 Eric Green, NHGRI\*\*](#)

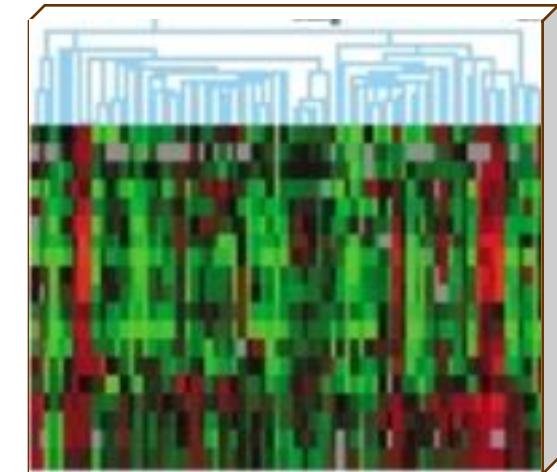


ATGTGCAATGCTT  
CGTTACGGCTCAA  
TATGCCGCAGTAA  
GCTGCAGTATCCG  
CCGCAGTAACCTGG  
GCCGCAG.....

Datos



Herramientas  
bioinformáticas



Conocimiento