



山東農業大學

SHANDONG AGRICULTURAL UNIVERSITY

生物信息学

第一章: 导论

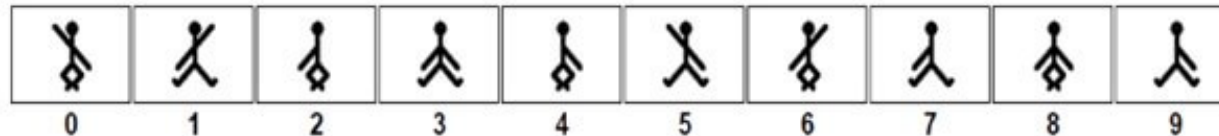
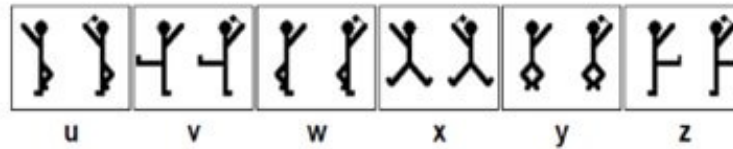
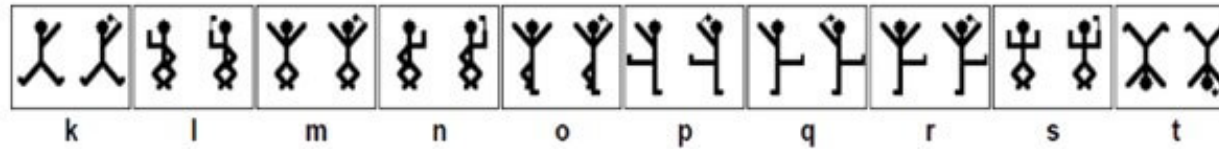
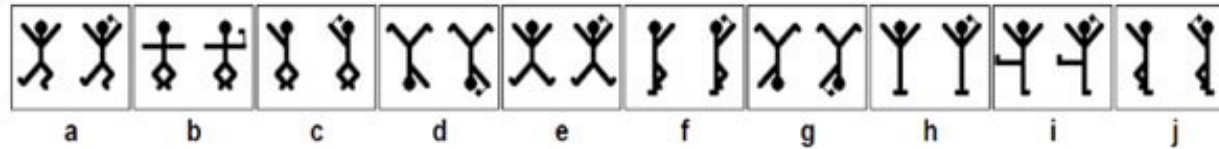
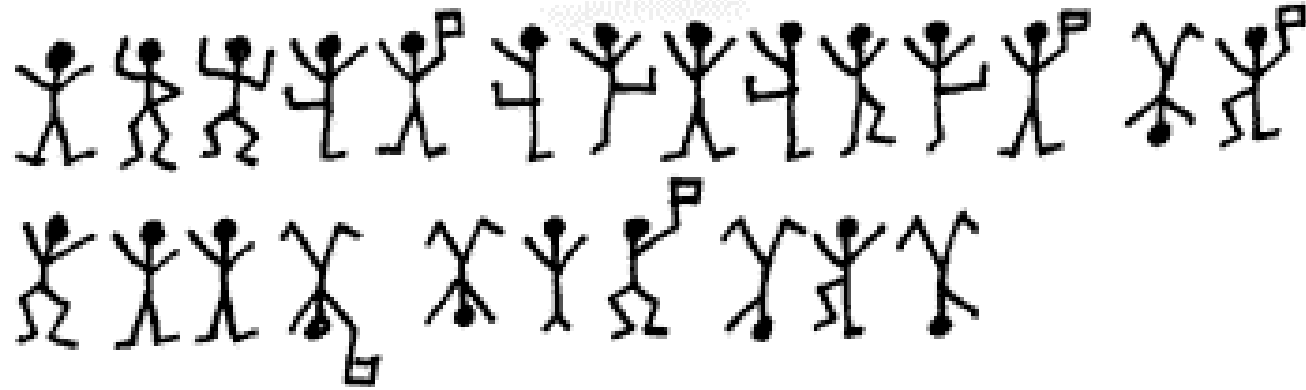
1. 生物信息学的概念

桂松涛

songtaogui@163.com

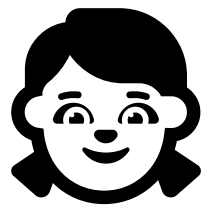


WHAT and WHY





WHAT and WHY



你猜我上午去药店的路上, 遇到谁了?



药店!药店!药店! 别问我怎么知道的



WHAT and WHY

A G C T T A C T A A T C C G G G C C G A A T T A G G T C
A G T T T A T T A A T T C G A G C T G A A C T A G G T C
A G T C T A T T A A T T C G A G C A G A A C T T G G T C
A G T T T A T T A A T T C G A G C T G A A C T T G G C C
A G T C T A C T A A T T C G A G C T G A A T T A G G T C
A G A T T A T T A A T T C G A G C T G A A C T T G G T C
A G A T T G C T A A T T C G A G C C G A A T T A G G T C
A G A T T A T T A A T C C G G G C T G A A T T A G G T C
A G T C T A T T A A T T C G A G C T G A A T T A G G A C
A G C T T A T T A A T T C G T G C T G A A C T C G G A C
A G C T T A T T A A T T C G A G C T G A A C T C G G A C
A G C T T A T T A A T T C G A G C C G A A C T C G G G C
A G T C T T T T A A T T C G A G C T G A A T T A G G A C

这段DNA序列是啥意思？

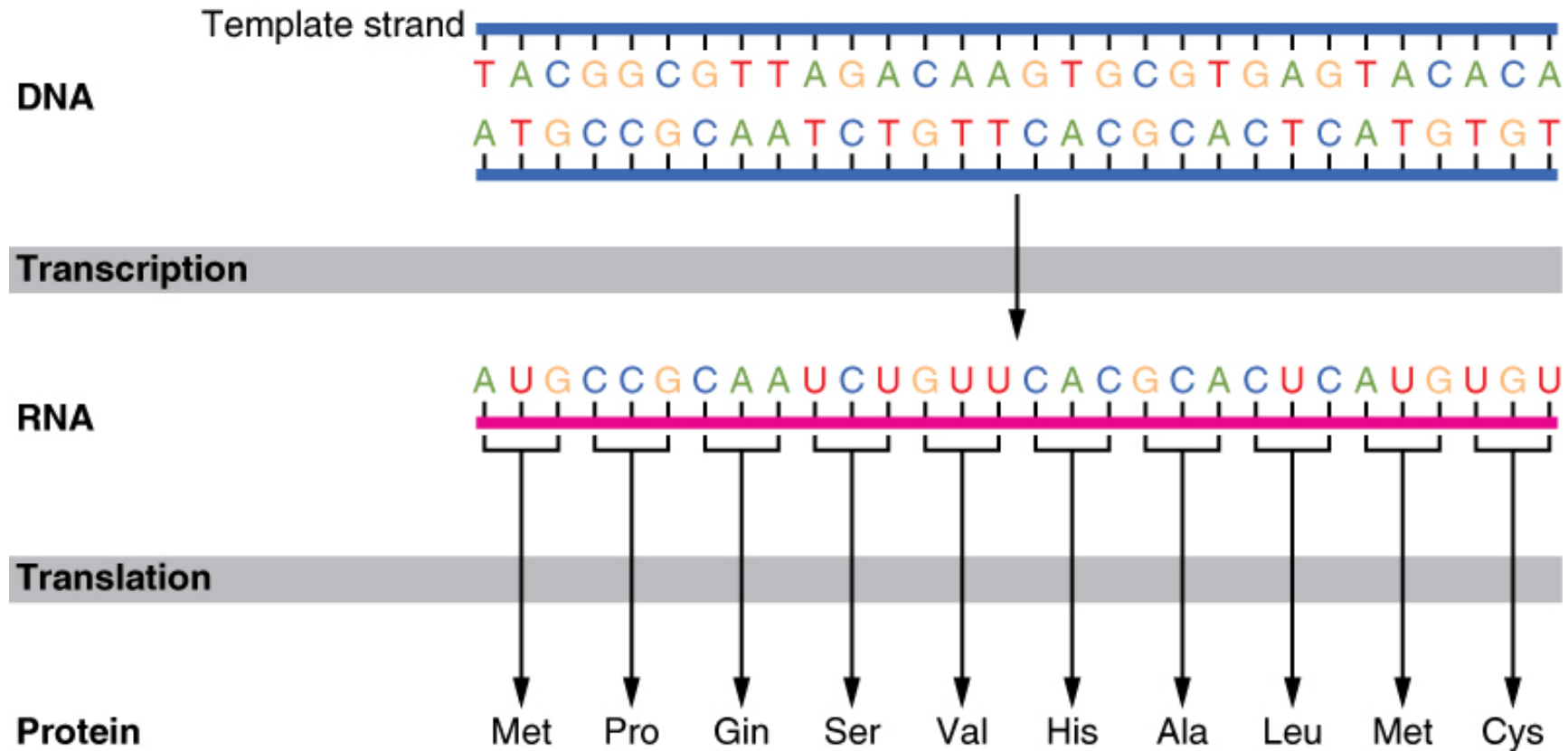
WHAT and WHY

tttctctttttttgtggatttttaaaggattttctacatagataaataaacacacataaacagttttacttcttt
cttttcaacctagactggatgcattttttgtttttgtttgtttgctttttaacttgctgcagtgactaga
gaatgtattgaagaatgattgttgaacaaaagcagtgagagtggacatccctgctttcccttgcttttaggg
ggaatgttttcagtccttccactatttaatatgatttttagctatgggtttatcctagatccctgttatcatgttg
aggaaattcccttctattttctagttttgttgagattttttaattcatgtgattgcgctatctggctttgctctca

- ◆ 人类基因组包含3G碱基对
- ◆ 两个人之间的基因组相似度为99.9%以上 (只有约3M碱基的差别)
- ◆ 平均每600 bp 就可能出现一个SNP
- ◆ 平均基因长度(含调控区) 27 Kb
- ◆ 基因编码区只占总长度的1.5%

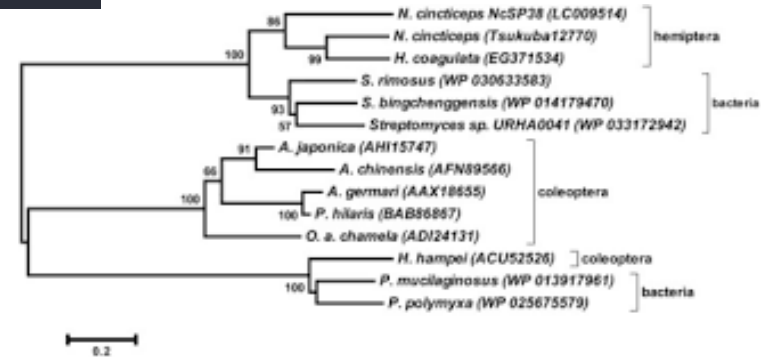
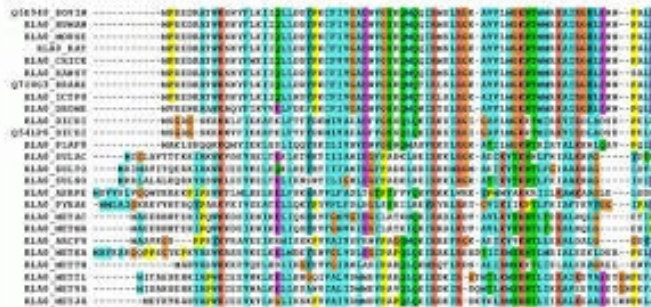
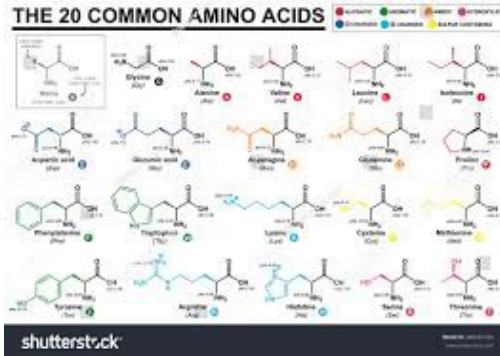
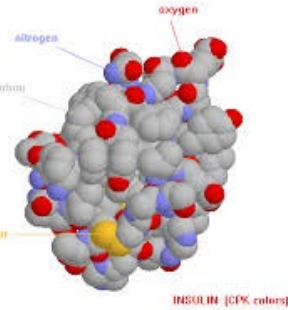
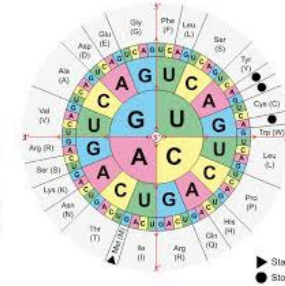
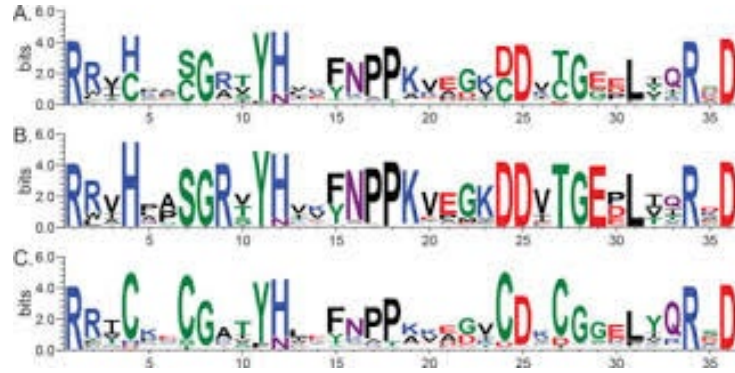


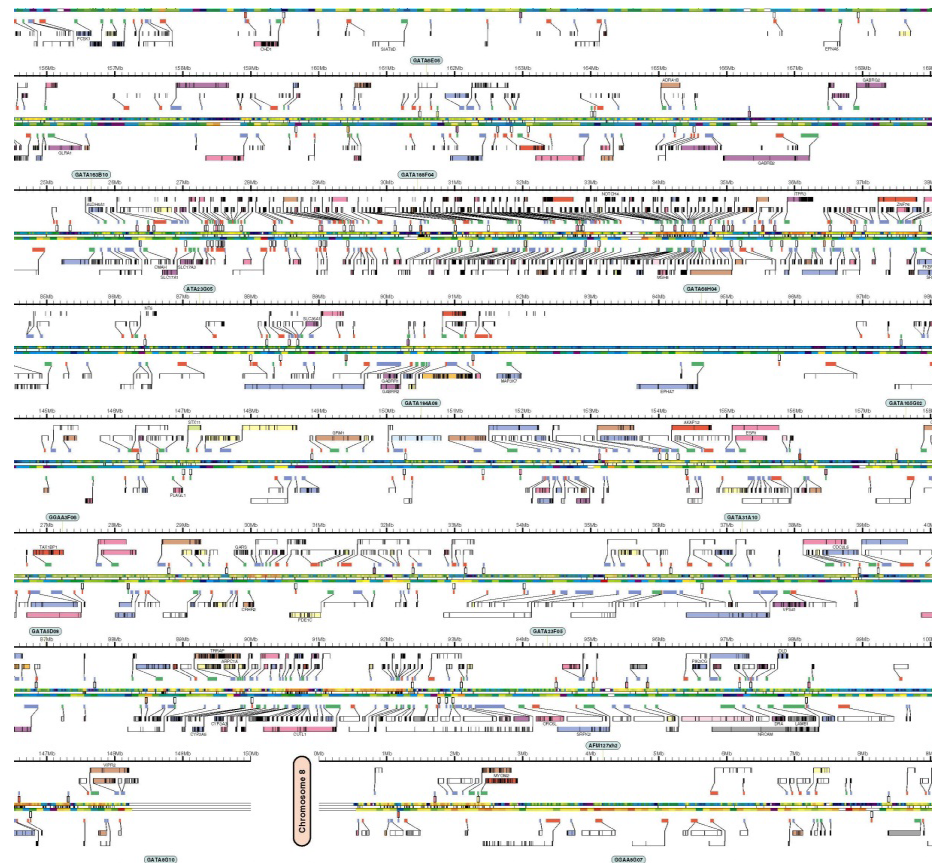
WHAT and WHY





WHAT and WHY

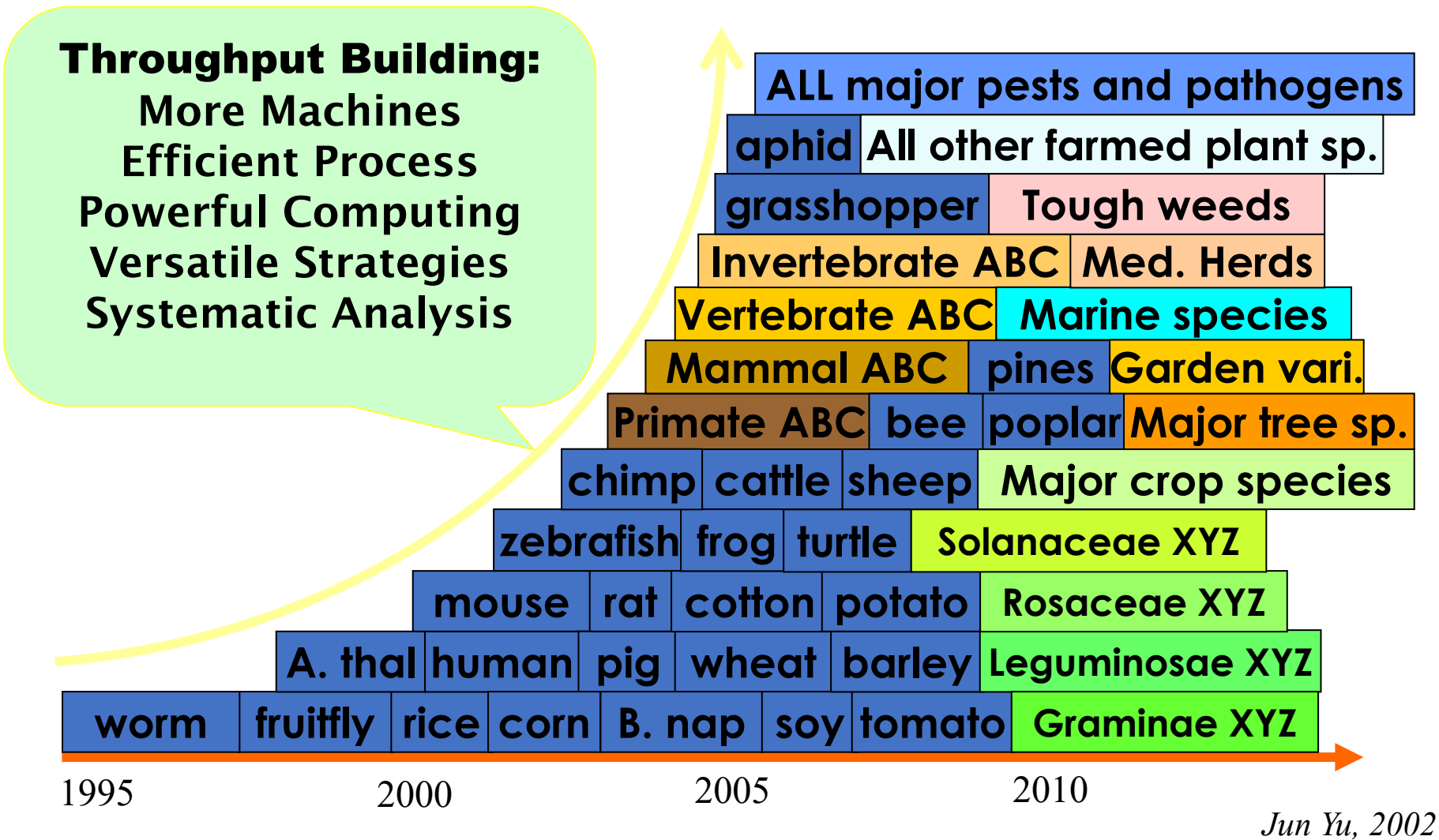




"The Sequence of the Human Genome"
Science, 2001



迅速增长的生物信息



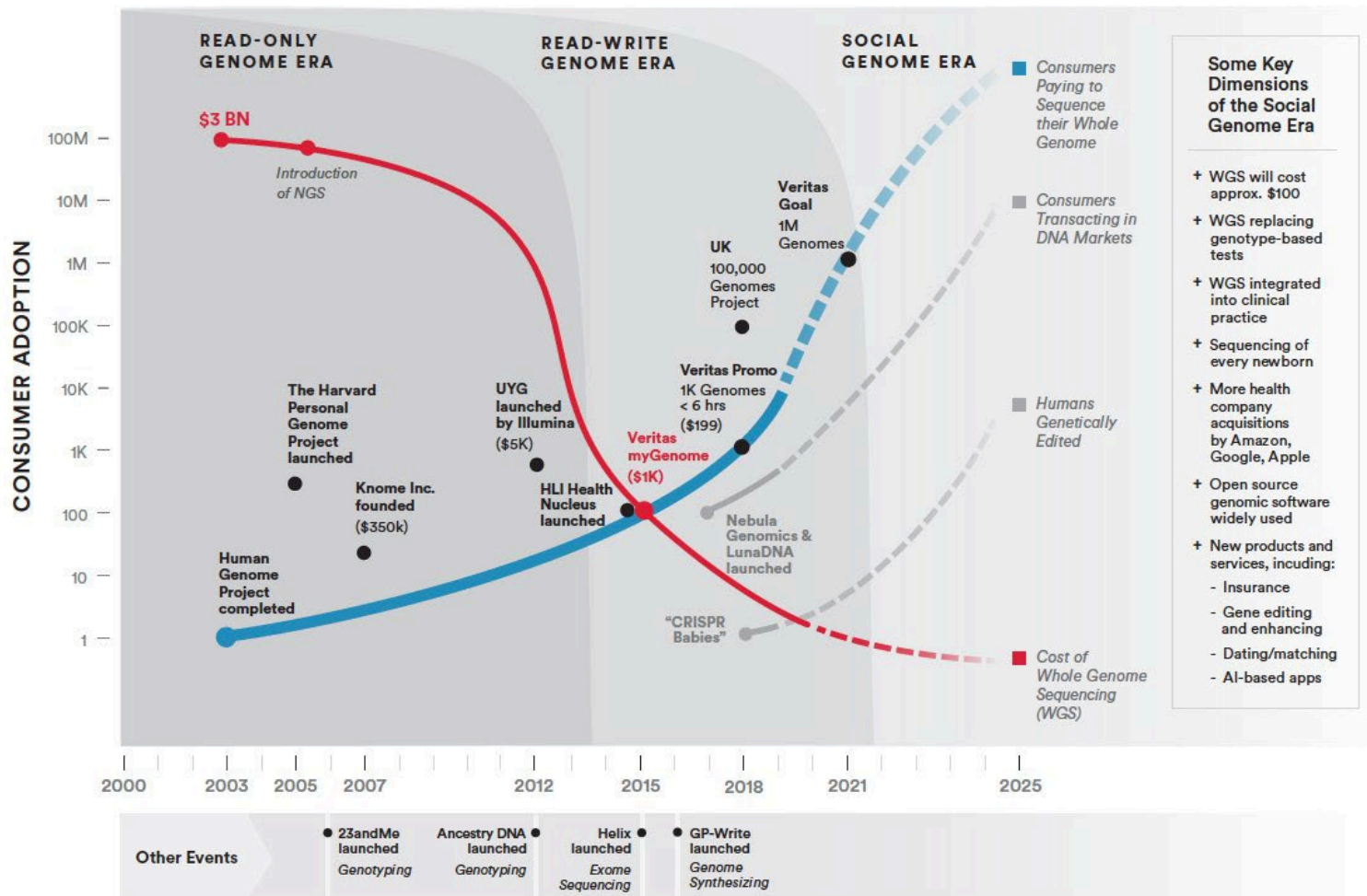


迅速增长的生物信息

The Emerging “Social Genome” v1.0 Feb 2019

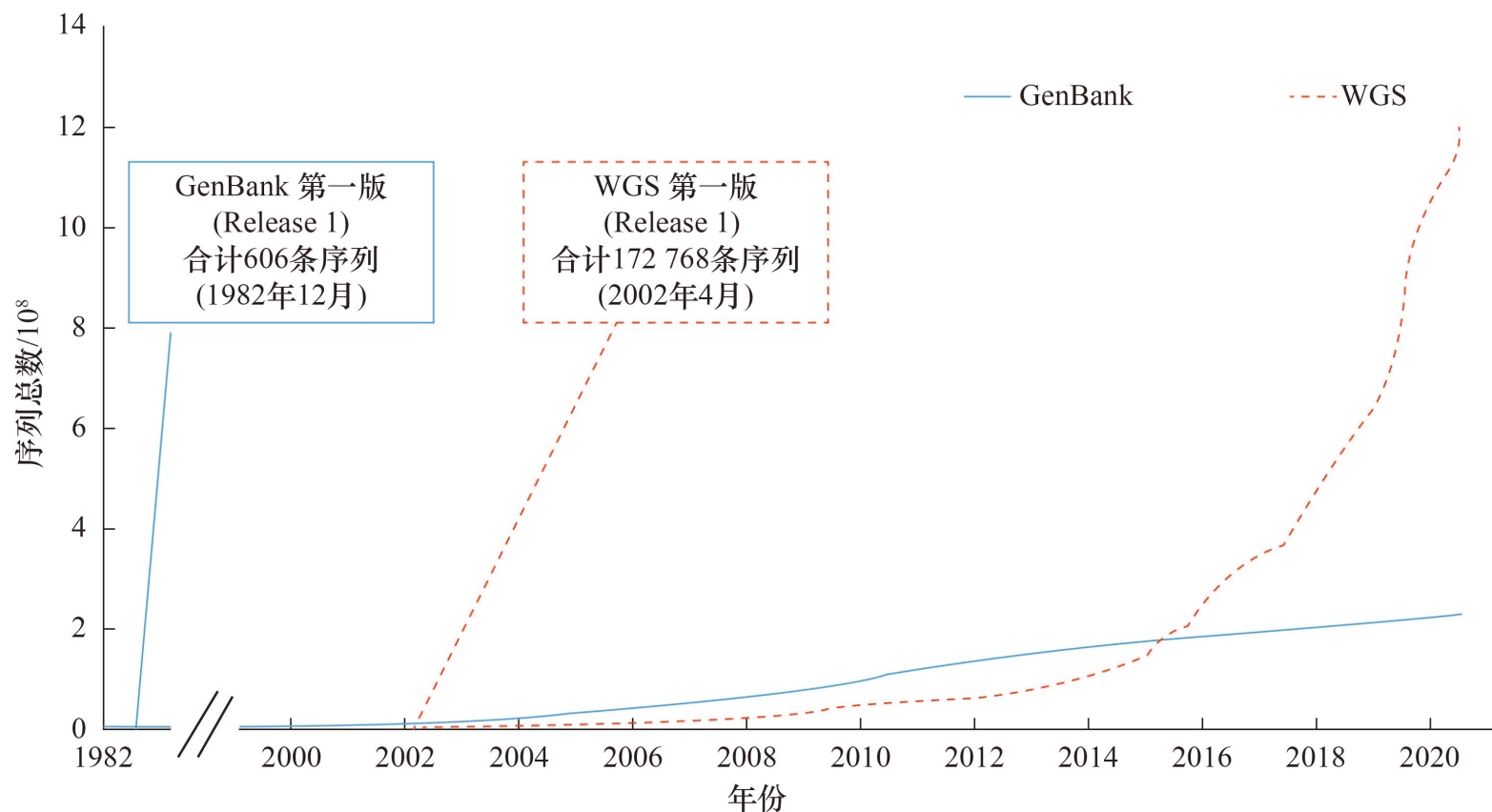
Illustrative Diagram

Some of these dates and figures are estimates based on publicly available information and conversations.





- 核酸数据库每14个月左右要翻一番
- 2024年, 公共核酸数据库: 8.2万亿个碱基对





- 高效的信息获取能力
- 批判性思维
- 自我学习和持续学习能力
- 保持好奇心和探索精神

信息爆炸时代如何保持洞察力？



- 有效管理
- 准确解读
- 充分使用



太长不看系列

“生物信息学是研究生物信息的采集、处理、存储、传播、分析和解释等的一门学科, 通过综合利用分子生物学、遗传学、计算机科学与技术, 来揭示大量且复杂的生物数据中蕴含的生物学奥秘。”

“生物信息学是利用信息来理解生物学的一门科学。”



山東農業大學

SHANDONG AGRICULTURAL UNIVERSITY

生物信息学

第一章: 导论

2. 生物信息学发展史

桂松涛

songtaogui@163.com



生物信息学名词的提出



宝琳·霍格维
(Paulien Hogeweg)
(1943-)

荷兰乌得勒支大学教授

- 1978年提出Bioinformatics一词(官方认证的最早)



林华安 博士
(Hwa A. Lim)
(1957-)

马来西亚华裔, 佛州超算中心遗传学与生物物理学部门主任

- 1990年组织了世界第一个国际生物学信息学学术会议, 催生了`生物信息学`一词的出现。



Bioinformatics

- Biology
- Information
- Mathematics



No. 4356 April 25, 1953

NATURE

737

equipment, and to Dr. G. E. R. Deacon and the captain and officers of R.R.S. *Discovery II* for their part in making the observations.

- ¹ Young, F. B., Gerrard, H., and Jevons, W., *Phil. Mag.*, **40**, 149 (1920).
² Longuet-Higgins, M. S., *Mon. Not. Roy. Astro. Soc., Geophys. Supp.*, **8**, 255 (1949).
³ Von Arx, W. S., Woods Hole Papers in Phys. Oceanogr. Meteor., **11** (3) (1950).
⁴ Ekman, V. W., *Arkiv. Mat. Astron. Fysik. (Stockholm)*, **2** (11) (1905).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three helical

is a residue on each chain every 3.4 Å. in the z-direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 10 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical z-co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 1; purine position 6 to pyrimidine position 6.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms

738

NATURE

April 25, 1953 VOL. 171

King's College, London. One of us (J.D.W.) has been aided by a fellowship from the National Foundation for Infantile Paralysis.

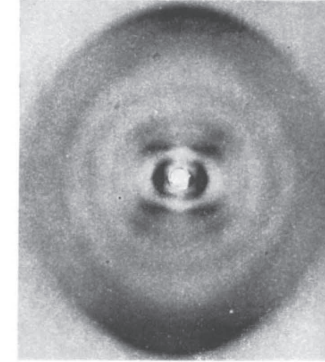
J. D. WATSON
F. H. C. CRICK

Medical Research Council Unit for the Study of the Molecular Structure of Biological Systems, Cavendish Laboratory, Cambridge, April 2.

- ¹ Pauling, L., and Corey, R. B., *Nature*, **171**, 346 (1953); *Proc. U.S. Nat. Acad. Sci.*, **38**, 81 (1953).
² Furlberg, S., *Acta Chem. Scand.*, **6**, 834 (1952).
³ Chargaff, E., for references see Zamenhof, S., Braverman, G., and Chargaff, E., *Biochim. et Biophys. Acta*, **9**, 402 (1952).
⁴ Wyatt, G. R., *J. Gen. Physiol.*, **36**, 201 (1952).
⁵ Astbury, W. T., *Symp. Soc. Exp. Biol.*, **1**, Nucleic Acid, 66 (Camb. Univ. Press, 1947).
⁶ Wilkins, M. H. F., and Randall, J. T., *Biochim. et Biophys. Acta*, **10**, 192 (1953).

Molecular Structure of Deoxypentose Nucleic Acids

WHILE the biological properties of deoxypentose



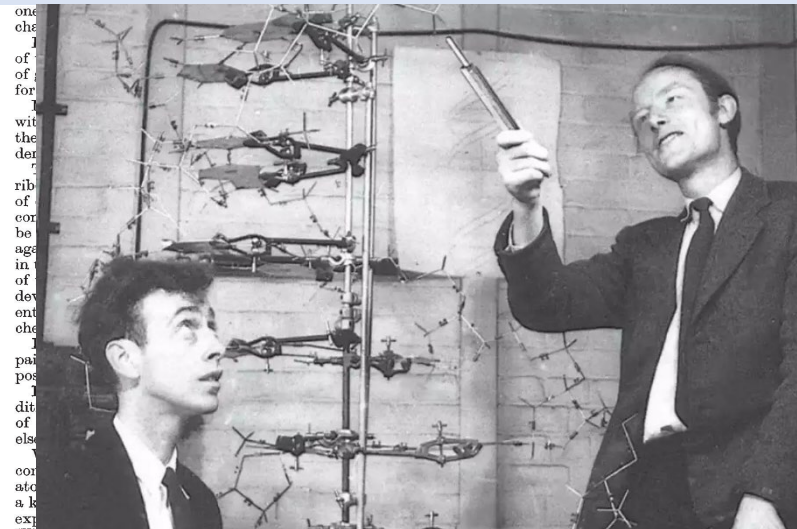
1953年，由沃森和克里克提出DNA双螺旋结构模型，开启了分子生物学时代

phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure as described is rather ill-defined, and for this reason we shall not comment on it.



This figure is purely diagrammatic. The two ribbons symbolize the two phosphate-sugar chains, and the horizontal rods the pairs of bases holding the chains together. The vertical line marks the fibre axis

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining β-D-deoxy-ribofuranose residues with 3',5' linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Furlberg's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Furlberg's 'standard configuration', the sugar being roughly perpendicular to the attached base. There



Wilkins, Dr. R. E. Franklin and their co-workers at

a group of parallel crystalline X-ray patterns, one of the longer sequence of not made

leic acid tion by as shown that the e inter- ~ 34 Å. eat of a ain con- cation as han the s on or helical

ed) that pattern helix is uniform spacing sity dis- portional function.

We will now briefly analyse in physical terms some of the effects of the shape and size of the repeat unit or nucleotide on the diffraction pattern. First, if the nucleotide consists of a unit having circular symmetry about an axis parallel to the helix axis, the whole diffraction pattern is modified by the form factor of the nucleotide. Second, if the nucleotide consists of a series of points on a radius at right-angles to the helix axis, the phases of radiation scattered by the helices of different diameter passing through each point are the same. Summation of the corresponding Bessel functions gives reinforcement for the inner-

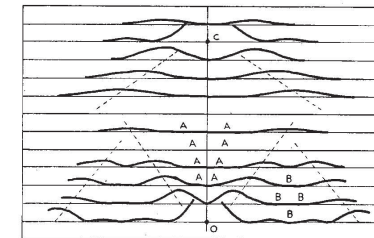


Fig. 2. Diffraction pattern of system of helices corresponding to structure of deoxypentose nucleic acid. The squares of Bessel functions are plotted about 0 on the equator and on the first, second, third and fifth layer lines for half of the nucleotide mass at 20 Å. diameter and remainder distributed along a radius, the mass at a given radius being proportional to the radius. About 0 on the tenth layer line similar functions are plotted for an outer diameter of 12 Å.

A straight line may be drawn approximately through



The Amide Groups of Insulin

By F. SANGER,* E. O. P. THOMPSON† AND RUTH KITAI
Department of Biochemistry, University of Cambridge

(Received 6 September 1954)

1955年, Sanger用二硝基氟苯 (FDNB) 法, 首次成功地完成了第一个蛋白质-牛胰岛素的序列分析

† Present address: Wool Textile Institute,
343 Royal Parade Parkville, Victoria, Australia

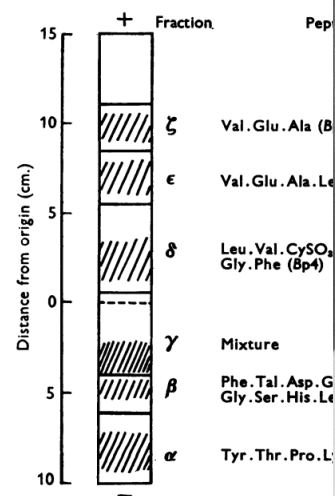
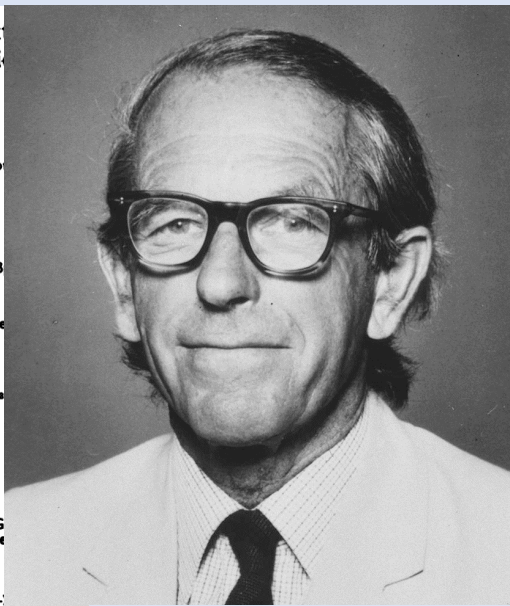
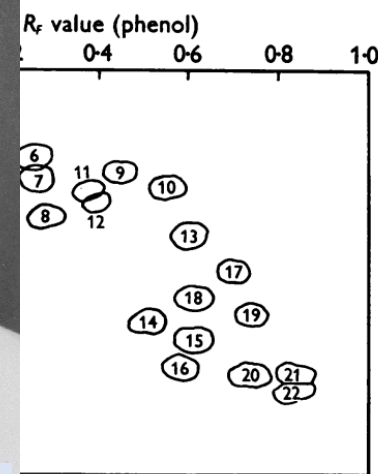


Fig. 1. Ionophoresis of peptic hydrolysate of insulin in 0.05 M ammonium acetate. 2.25 V, 20 min.



Frederick Sanger

hydrolyde Chibnall & Rees (1952)
the three aspartic acids, one of which is



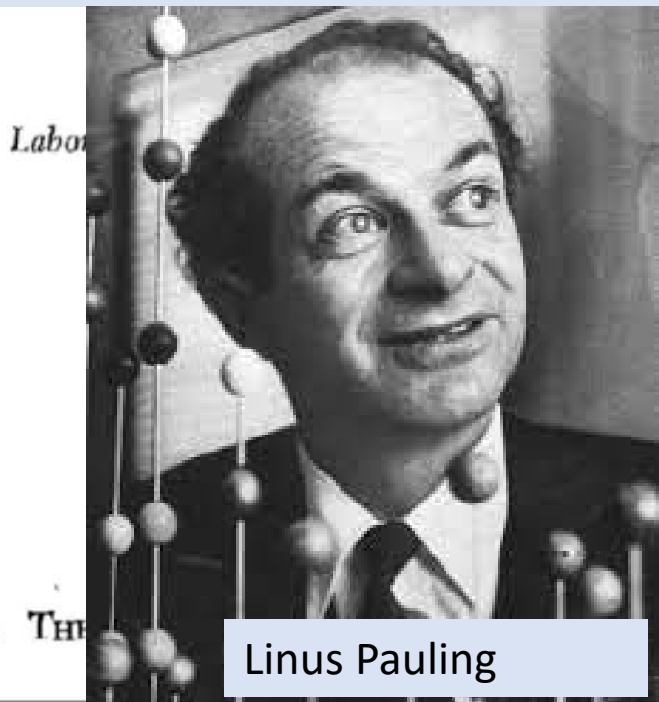
of mould protease hydrolysate of insulin (Expt. 4m) (see Table 1).



Evolving genes and proteins, 1965

Evolutionary Divergence and Convergence in Proteins

1965年，祖卡坎德尔和鲍林提出的“分子钟”理论



Linus Pauling



1966年，我国第一次人工合成了胰岛素





Proc. Natl. Acad. Sci. USA
Vol. 74, No. 12, pp. 5463–5467, December 1977
Biochemistry

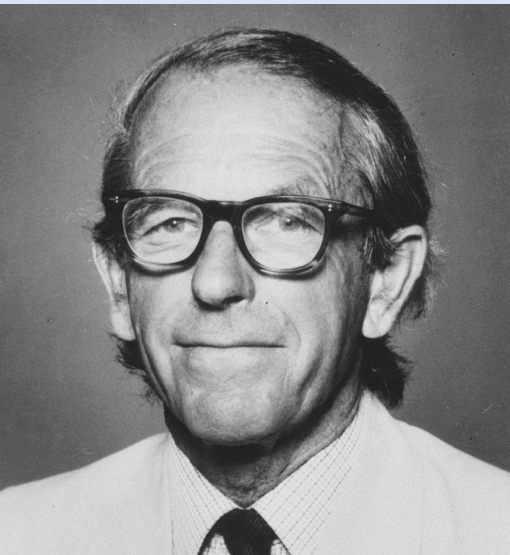
1977年，桑格等发表双脱氧链末端终止法，测定 ϕ X174序列。

Medical Research Council Laboratory of Molecular Biology, Cambridge

Contributed by F. Sanger, October 3, 1977

ABSTRACT A new method for determining the sequence of DNA is described. It is a "plus and minus" method [Sanger, F. & Coulson, 1975, *Nature*, **251**, 441–448] but makes use of the 2',3'-dideoxynucleoside analogues of the normal deoxynucleosides which act as specific chain-terminating agents for DNA polymerase. The technique has been applied to the sequencing of bacteriophage ϕ X174 and is more rapid than either the plus or the minus method.

The "plus and minus" method (1) is a simple technique that has made possible the sequencing of the genome of bacteriophage ϕ X174. It depends on the use of DNA polymerase and the synthesis of complementary regions of the DNA under controlled conditions. This method is considerably more rapid than the available techniques, neither the "plus" nor the "minus" method is completely accurate, and in some cases both must be used together, and in some cases sequence data are necessary. W. M. Barnes has recently developed a third method, involving the substitution, which has certain advantages over the plus and minus



Frederick Sanger, again!

of ribose in which the 3'-hydroxyl group is oriented in the opposite position with respect to the 2'-hydroxyl group. The 2',3'-dideoxynucleoside (ara) nucleotides act as chain terminating inhibitors for *Escherichia coli* DNA polymerase I in a manner similar to that of ddT (4), although synthesized chains ending in ddT are further extended by some mammalian DNA polymerases (5). In order to obtain a suitable pattern of bands for sequencing an extensive sequence can be read it is necessary to use a ratio of terminating triphosphate to normal triphosphate of only partial incorporation of the terminator. In the case of the dideoxy derivatives this ratio is about 100, and in the case of the arabinoside derivatives about 5000.

METHODS

Preparation of the Triphosphate Analogues. The preparation of the triphosphate analogues has been described (6, 7), and the material is commercially available. ddA has been prepared by the method of Tener (8). We essentially followed their procedure and used the methods of Tener (9) and of Hoard and Ott (10) to convert it to the triphosphate, which was then purified on



1988年，人类基因组计划提出

Perspective

A Turning Point in Cancer Research: Sequencing the Human Genome

RENATO DULBECCO

ONE OF THE GOALS OF CANCER RESEARCH IS TO ASCERTAIN the mechanisms of cancer. Efforts in this direction have been made by using model systems of limited complexity, such as cancer cells in vitro and oncogenic viruses. The use of cell cultures avoided the complexity of the whole animal but not the complexity of the animal genome. The use of oncogenic viruses seemed to circumvent this complexity by replacing it with the extraordinary simplicity of the viral genome. This simplicity made the study of viruses very productive. The persistence of the transformed state in a cell clone could be explained by the persistence of the viral genome in cells (1); genetic and molecular results showed that transformation is the consequence of the expression of one or a few viral genes. Finally, the viral transforming genes, or "oncogenes," and the proteins they specify were identified. The crowning development was the demonstration that in retroviruses the oncogenes are picked up from the cellular genome during the viruses' most recent history (2). As a result of these studies, cancer seemed to be locked to the expression of some viral gene; the possibility of a "hit-and-run" mechanism, in which the virus alters the cell and then vanishes, seemed excluded. Two types of oncogenes were identified: some which immortalize cells, and others which make them tumorigenic (3). In most cases oncogenes of both types are needed to cause a continuously growing tumor.

Subsequent work, however, blurred the distinction between immortalizing and transforming oncogenes by showing that their effects differ in primary cultures or permanent lines and in cells of different species (4). These findings suggested that the state of the cellular genes is important for the effect of oncogenes, in agreement with the great differences in cancer incidence and in the effects of chemical or viral carcinogens in different species.

These studies dealt with the initial cancer events. But natural cancers evolve slowly toward malignancy through many definable stages in a process called "progression" (5), which is the least understood but probably the most crucial phase in the generation of malignancy. Progression generates the marked heterogeneity of cancers (6) and their many chromosomal abnormalities (7); it must be differentiated from the initial action of oncogenes (8). Progression is observed in cells transformed by viruses. This is the case, for instance, of bursal lymphomas induced by avian leukosis viruses (9), of viral T-cell lymphomas in mice (10), and of leukemogenesis by Friend leukemia virus in cultures of mouse bone marrow cells (11). Stepwise transformation is observed also with DNA viruses (12). Fibroblastic cells from a variety of organs of a transgenic mouse

containing *myc* and simian virus 40 (SV40) sequences, although expressing SV40 T antigen, were normal but became gradually transformed upon cultivation (13). In all these cases cellular changes occurring during culture growth determined full transformation. The "hit-and-run" hypothesis of viral transformation must be reconsidered.

A clue as to what these changes are is obtained by examining the heterogeneity of chemically induced rat mammary carcinomas with respect to several well-characterized markers. The expression of the markers is altered in different ways in different parts of the same cancer; the alterations seem to be clonal, being uniform in small parts of a tumor but different in adjacent parts (14). The closeness of the parts makes it unlikely that the differences are due to the environment; it is more likely that they are caused by structural changes of the genes, as is also suggested by the chromosomal rearrangements observed in cancers (15) and by the finding that each chemically or radiation-induced mouse sarcoma expresses a different class I major histocompatibility antigen, probably produced through gene rearrangement (16).

A major gap in our understanding of cancer is how the activity of an oncogene is related to the events of progression. But the first task is to ascertain whether the DNA of an advanced cancer is as heterogeneous as the phenotype of its cells. If it is so, a new field of cancer research opens up, possibly leading to the discovery of the genes whose activity or inactivity is responsible for infiltration and metastasis.

We are at a turning point in the study of tumor virology and cancer in general. If we wish to learn more about cancer, we must now concentrate on the cellular genome. We are back to where cancer research started, but the situation is drastically different because we have new knowledge and crucial tools, such as DNA cloning. We have two options: either to try to discover the genes important in malignancy by a piecemeal approach, or to sequence the whole genome of a selected animal species. The former approach seems less formidable, but it will still require a vast investment of research, especially if the important genes differ in cancers of different organs and if they encode regulatory proteins. A major difficulty for conventional approaches is the heterogeneity of tumors and the lack of cultures representative of the various cell types present in a cancer. I think that it will be far more useful to begin by sequencing the cellular genome. The sequence will make it possible to prepare probes for all the genes and to classify them for their expression in various cell types at the level of individual cells by means of cytological hybridization. The classification of the genes will facilitate the identification of those involved in progression.

In which species should this effort be made? If we wish to understand human cancer, it should be made in humans because the genetic control of cancer seems to be different in different species. Research on human cancer would receive a major boost from the detailed knowledge of DNA. Humans would become the preferred experimental species for cancer research with cells in culture or in immunodeficient mice. Because cancer could be defined in molecular terms, the agents capable of inducing cancer in humans could be identified by the combination of in vitro and epidemiological studies. Knowledge of the genes involved in progression would open new therapeutic approaches, which might lead to a general cancer cure if progression has common features in all cancers.

Knowledge of the genome and availability of probes for any gene would also be crucial for progress in human physiology and pathology outside cancer; for instance, for learning about the regulation of individual genes in various cell types. Many fields of

The author is in the Monoclonal Antibody Laboratory of the Armand Hammer Cancer Center, the Salk Institute, La Jolla, CA 92037.

Downloaded from www.sciencemag.org on March 11, 2014

人类基因组计划

开始于1988年

由美国能源部和国家医学研究院发起

美国、英国、法国、德国、日本和中国参加

经费三十亿美元

主要目标：

1、人类基因组测序

2、各种模式生物基因组测序

3、推动全基因组水平的高通量技术的发展



RESEARCH ARTICLE

Whole-Genome Random Sequencing and Assembly of *Haemophilus influenzae* Rd

Robert D. Fleischmann, Mark D. Adams, Owen White, Rebecca A. Clayton, Ewen F. Kirkness, Anthony R. Kerlavage, Carol J. Bult, Jean-Francois Tomb, Brian A. Dougherty, Joseph M. Muzny, Keith McKernan, Gregory Sutton

1995年, *H. influenza* (流感嗜血杆菌)基因组: 第一个测序成功的基因组。

Hamilton O. Smith, J. Craig Venter†

An approach for genome analysis based on sequencing and assembly of unselected pieces of DNA from the whole chromosome has been applied to obtain the complete nucleotide sequence (1,830,137 base pairs) of the genome from the bacterium *Haemophilus influenzae* Rd. This approach eliminates the need for initial mapping efforts and is therefore applicable to the vast array of microbial species for which genome maps are unavailable. The *H. influenzae* Rd genome sequence (Genome Sequence DataBase accession number L42023) represents the only complete genome sequence from a free-living organism.

A prerequisite to understanding the complete biology of an organism is the determination of its entire genome sequence. Several viral and organellar genomes have been sequenced, including those of *Bacteriophage lambda* (1), *Adenovirus* (2), *Herpesvirus* (3), *Poliovirus* (4), *Rotavirus* (5), *Simian virus 40* (6), *Human immunodeficiency virus* (7), *Human papillomavirus* (8), *Human cytomegalovirus* (9), *Human herpesvirus 8* (10), and *Homo sapiens* (11). These projects, as well as viral genome sequencing, have been based primarily on the sequencing of clones usually derived from extensively mapped restriction fragments (12). The first complete genome sequence of a free-living organism was that of *Haemophilus influenzae* Rd (13). This bacterium is a natural host is human. Six *H. influenzae* serotype strains (a through f) have been identified on the basis of immunologically distinct capsular polysaccharide antigens. Non-typeable strains also exist and are distinguished by their lack of detectable capsular polysaccharide. They are commensal residents of the upper respiratory mucosa of children and adults and cause otitis media and respiratory tract infections, mostly in

ly reduced the incidence of the disease in Europe and North America.

Genome sequencing. The strategy for a shotgun approach to whole genome sequencing is outlined in Table 1. The theory follows from the Lander and Waterman (14) application of the equation for the Poisson distribution. The probability that a base is not sequenced is $P_0 = e^{-m}$, where m is the sequence coverage. Thus after 1.83 Mb of sequence has been randomly generated for the *H. influenzae* genome ($m = 1, 1 \times$ coverage), $P_0 = e^{-1} = 0.37$ and approximately 37 percent of the genome is unsequenced. Fivefold coverage (approximately 9500 clones sequenced from both insert ends and an average sequence read length of 460 bp) will give $P_0 = e^{-5} = 0.0067$, or approximately 0.7 percent of the genome is unsequenced.

ly reduced the incidence of the disease in Europe and North America.

articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

The rediscovery of Mendel's laws of heredity in the open 20th century¹⁻³ sparked a scientific quest to understand the nature and content of genetic information that has been passed on from one generation to the next. The scientific pursuit falls naturally into four main phases, corresponding to the four quarters of the century. The first established the cellular basis of heredity: the chromosomes. The second defined the molecular basis of heredity: the DNA double helix. The third unlocked the informational basis of heredity, with the discovery of the biological mechanism by which cells read the information contained in genes and with the invention of the recombinant DNA technologies of cloning and sequencing by which scientists can do the same.

The last quarter of a century has been marked by a relentless drive to decipher first genes and then entire genomes, spawning the field of genomics. The fruits of this work already include the genome sequences of 599 viruses and viroids, 205 naturally occurring plasmids, 185 organelles, 31 eubacteria, seven archaea, one fungus, two animals and one plant.

Here we report the results of a collaboration involving 20 groups from the United States, the United Kingdom, Japan, France, Germany and China to produce a draft sequence of the human genome. The draft genome sequence was generated from a physical map covering more than 96% of the euchromatic part of the human genome and, together with additional sequence in public databases, it covers about 94% of the human genome. The sequence was produced over a relatively short period, with coverage rising from about 10% to more than 90% over roughly fifteen months. The sequence data have been made available without restriction and updated daily throughout the project. The task ahead is to produce a finished sequence, by closing all gaps and resolving all ambiguities. Already about one billion bases are in final form and the task of bringing the vast majority of the sequence to this standard is now straightforward and should proceed rapidly.

The sequence of the human genome is of interest in several respects. It is the largest genome to be extensively sequenced so far, being 25 times as large as any previously sequenced genome and eight times as large as the sum of all such genomes. It is the first vertebrate genome to be extensively sequenced. And, uniquely, it is the genome of our own species.

Much work remains to be done to produce a complete finished sequence, but the vast trove of information that has become available through this collaborative effort allows a global perspective on the human genome. Although the details will change as the sequence is finished, many points are already clear.

- The genomic landscape shows marked variation in the distribution of a number of features, including genes, transposable elements, GC content, CpG islands and recombination rate. This gives us important clues about function. For example, the developmentally important HOX gene clusters are the most repeat-poor regions of the human genome, probably reflecting the very complex

2001年, 人类基因组草图公布。

The human genome is more complex than those of invertebrates. This is due in part to the presence of vertebrate-specific protein domains and motifs (an estimated 7% of the total), but more to the fact that vertebrates appear to have arranged pre-existing components into a richer collection of domain architectures.

- Hundreds of human genes appear likely to have resulted from horizontal transfer from bacteria at some point in the vertebrate lineage. Dozens of genes appear to have been derived from transposable elements.

- Although about half of the human genome derives from transposable elements, there has been a marked decline in the overall activity of such elements in the hominid lineage. DNA transposons appear to have become completely inactive and long-terminal repeat (LTR) retrotransposons may also have done so.

- The pericentromeric and subtelomeric regions of chromosomes are filled with large recent segmental duplications of sequence from elsewhere in the genome. Segmental duplication is much more frequent in humans than in yeast, fly or worm.

- Analysis of the organization of Alu elements explains the long-standing mystery of their surprising genomic distribution, and suggests that there may be strong selection in favour of preferential retention of Alu elements in GC-rich regions and that these 'selfish' elements may benefit their human hosts.

- The mutation rate is about twice as high in male as in female meiosis, showing that most mutation occurs in males.

- Cytogenetic analysis of the sequenced clones confirms suggestions that large GC-poor regions are strongly correlated with 'dark G-bands' in karyotypes.

- Recombination rates tend to be much higher in distal regions (around 20 megabases (Mb)) of chromosomes and on shorter chromosome arms in general, in a pattern that promotes the occurrence of at least one crossover per chromosome arm in each meiosis.

- More than 1.4 million single nucleotide polymorphisms (SNPs) in the human genome have been identified. This collection should allow the initiation of genome-wide linkage disequilibrium mapping of the genes in the human population.

In this paper, we start by presenting background information on the project and describing the generation, assembly and evaluation of the draft genome sequence. We then focus on an initial analysis of the sequence itself: the broad chromosomal landscape; the repeat elements and the rich palaeontological record of evolutionary and biological processes that they provide; the human genes and proteins and their differences and similarities with those of other

THE HUMAN GENOME

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹ Peter W. Li,¹ Richard J. Mural,¹ Granger G. Sutton,¹ Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹ Robert A. Holt,¹ Jeannine D. Gocayne,¹ Peter Amanatides,¹ Richard M. Ballew,¹ Daniel H. Huson,¹ Jennifer Russo Wortman,¹ Qing Zhang,¹ Chinnappa D. Kodira,¹ Xiangqun H. Zheng,¹ Lin Chen,¹ Marian Skupski,¹ Gangadharan Subramanian,¹ Paul D. Thomas,¹ Jinghui Zhang,¹ George L. Gabor Miklos,² Catherine Nelson,³ Samuel Broder,¹ Andrew G. Clark,⁴ Joe Nadeau,⁵ Victor A. McKusick,⁶ Norton Zinder,⁷ Arnold J. Levine,⁷ Richard J. Roberts,⁸ Mel Simon,⁹ David H. Ledbetter,¹⁰ Robert H. Waterman,¹¹ David L. Page,¹² David R. Bentley,¹³ Michael R. Brown,¹⁴ Michael J. Morgan,¹⁵ Michael J. Smith,¹⁶ Michael J. Morgan,¹⁷ Michael J. Morgan,¹⁸ Michael J. Morgan,¹⁹ Michael J. Morgan,²⁰ Michael J. Morgan,²¹ Michael J. Morgan,²² Michael J. Morgan,²³ Michael J. Morgan,²⁴ Michael J. Morgan,²⁵ Michael J. Morgan,²⁶ Michael J. Morgan,²⁷ Michael J. Morgan,²⁸ Michael J. Morgan,²⁹ Michael J. Morgan,³⁰ Michael J. Morgan,³¹ Michael J. Morgan,³² Michael J. Morgan,³³ Michael J. Morgan,³⁴ Michael J. Morgan,³⁵ Michael J. Morgan,³⁶ Michael J. Morgan,³⁷ Michael J. Morgan,³⁸ Michael J. Morgan,³⁹ Michael J. Morgan,⁴⁰ Michael J. Morgan,⁴¹ Michael J. Morgan,⁴² Michael J. Morgan,⁴³ Michael J. Morgan,⁴⁴ Michael J. Morgan,⁴⁵ Michael J. Morgan,⁴⁶ Michael J. Morgan,⁴⁷ Michael J. Morgan,⁴⁸ Michael J. Morgan,⁴⁹ Michael J. Morgan,⁵⁰ Michael J. Morgan,⁵¹ Michael J. Morgan,⁵² Michael J. Morgan,⁵³ Michael J. Morgan,⁵⁴ Michael J. Morgan,⁵⁵ Michael J. Morgan,⁵⁶ Michael J. Morgan,⁵⁷ Michael J. Morgan,⁵⁸ Michael J. Morgan,⁵⁹ Michael J. Morgan,⁶⁰ Michael J. Morgan,⁶¹ Michael J. Morgan,⁶² Michael J. Morgan,⁶³ Michael J. Morgan,⁶⁴ Michael J. Morgan,⁶⁵ Michael J. Morgan,⁶⁶ Michael J. Morgan,⁶⁷ Michael J. Morgan,⁶⁸ Michael J. Morgan,⁶⁹ Michael J. Morgan,⁷⁰ Michael J. Morgan,⁷¹ Michael J. Morgan,⁷² Michael J. Morgan,⁷³ Michael J. Morgan,⁷⁴ Michael J. Morgan,⁷⁵ Michael J. Morgan,⁷⁶ Michael J. Morgan,⁷⁷ Michael J. Morgan,⁷⁸ Michael J. Morgan,⁷⁹ Michael J. Morgan,⁸⁰ Michael J. Morgan,⁸¹ Michael J. Morgan,⁸² Michael J. Morgan,⁸³ Michael J. Morgan,⁸⁴ Michael J. Morgan,⁸⁵ Michael J. Morgan,⁸⁶ Michael J. Morgan,⁸⁷ Michael J. Morgan,⁸⁸ Michael J. Morgan,⁸⁹ Michael J. Morgan,⁹⁰ Michael J. Morgan,⁹¹ Michael J. Morgan,⁹² Michael J. Morgan,⁹³ Michael J. Morgan,⁹⁴ Michael J. Morgan,⁹⁵ Michael J. Morgan,⁹⁶ Michael J. Morgan,⁹⁷ Michael J. Morgan,⁹⁸ Michael J. Morgan,⁹⁹ Michael J. Morgan,¹⁰⁰ Michael J. Morgan,¹⁰¹ Michael J. Morgan,¹⁰² Michael J. Morgan,¹⁰³ Michael J. Morgan,¹⁰⁴ Michael J. Morgan,¹⁰⁵ Michael J. Morgan,¹⁰⁶ Michael J. Morgan,¹⁰⁷ Michael J. Morgan,¹⁰⁸ Michael J. Morgan,¹⁰⁹ Michael J. Morgan,¹¹⁰ Michael J. Morgan,¹¹¹ Michael J. Morgan,¹¹² Michael J. Morgan,¹¹³ Michael J. Morgan,¹¹⁴ Michael J. Morgan,¹¹⁵ Michael J. Morgan,¹¹⁶ Michael J. Morgan,¹¹⁷ Michael J. Morgan,¹¹⁸ Michael J. Morgan,¹¹⁹ Michael J. Morgan,¹²⁰ Michael J. Morgan,¹²¹ Michael J. Morgan,¹²² Michael J. Morgan,¹²³ Michael J. Morgan,¹²⁴ Michael J. Morgan,¹²⁵ Michael J. Morgan,¹²⁶ Michael J. Morgan,¹²⁷ Michael J. Morgan,¹²⁸ Michael J. Morgan,¹²⁹ Michael J. Morgan,¹³⁰ Michael J. Morgan,¹³¹ Michael J. Morgan,¹³² Michael J. Morgan,¹³³ Michael J. Morgan,¹³⁴ Michael J. Morgan,¹³⁵ Michael J. Morgan,¹³⁶ Michael J. Morgan,¹³⁷ Michael J. Morgan,¹³⁸ Michael J. Morgan,¹³⁹ Michael J. Morgan,¹⁴⁰ Michael J. Morgan,¹⁴¹ Michael J. Morgan,¹⁴² Michael J. Morgan,¹⁴³ Michael J. Morgan,¹⁴⁴ Michael J. Morgan,¹⁴⁵ Michael J. Morgan,¹⁴⁶ Michael J. Morgan,¹⁴⁷ Michael J. Morgan,¹⁴⁸ Michael J. Morgan,¹⁴⁹ Michael J. Morgan,¹⁵⁰ Michael J. Morgan,¹⁵¹ Michael J. Morgan,¹⁵² Michael J. Morgan,¹⁵³ Michael J. Morgan,¹⁵⁴ Michael J. Morgan,¹⁵⁵ Michael J. Morgan,¹⁵⁶ Michael J. Morgan,¹⁵⁷ Michael J. Morgan,¹⁵⁸ Michael J. Morgan,¹⁵⁹ Michael J. Morgan,¹⁶⁰ Michael J. Morgan,¹⁶¹ Michael J. Morgan,¹⁶² Michael J. Morgan,¹⁶³ Michael J. Morgan,¹⁶⁴ Michael J. Morgan,¹⁶⁵ Michael J. Morgan,¹⁶⁶ Michael J. Morgan,¹⁶⁷ Michael J. Morgan,¹⁶⁸ Michael J. Morgan,¹⁶⁹ Michael J. Morgan,¹⁷⁰ Michael J. Morgan,¹⁷¹ Michael J. Morgan,¹⁷² Michael J. Morgan,¹⁷³ Michael J. Morgan,¹⁷⁴ Michael J. Morgan,¹⁷⁵ Michael J. Morgan,¹⁷⁶ Michael J. Morgan,¹⁷⁷ Michael J. Morgan,¹⁷⁸ Michael J. Morgan,¹⁷⁹ Michael J. Morgan,¹⁸⁰ Michael J. Morgan,¹⁸¹ Michael J. Morgan,¹⁸² Michael J. Morgan,¹⁸³ Michael J. Morgan,¹⁸⁴ Michael J. Morgan,¹⁸⁵ Michael J. Morgan,¹⁸⁶ Michael J. Morgan,¹⁸⁷ Michael J. Morgan,¹⁸⁸ Michael J. Morgan,¹⁸⁹ Michael J. Morgan,¹⁹⁰ Michael J. Morgan,¹⁹¹ Michael J. Morgan,¹⁹² Michael J. Morgan,¹⁹³ Michael J. Morgan,¹⁹⁴ Michael J. Morgan,¹⁹⁵ Michael J. Morgan,¹⁹⁶ Michael J. Morgan,¹⁹⁷ Michael J. Morgan,¹⁹⁸ Michael J. Morgan,¹⁹⁹ Michael J. Morgan,²⁰⁰ Michael J. Morgan,²⁰¹ Michael J. Morgan,²⁰² Michael J. Morgan,²⁰³ Michael J. Morgan,²⁰⁴ Michael J. Morgan,²⁰⁵ Michael J. Morgan,²⁰⁶ Michael J. Morgan,²⁰⁷ Michael J. Morgan,²⁰⁸ Michael J. Morgan,²⁰⁹ Michael J. Morgan,²¹⁰ Michael J. Morgan,²¹¹ Michael J. Morgan,²¹² Michael J. Morgan,²¹³ Michael J. Morgan,²¹⁴ Michael J. Morgan,²¹⁵ Michael J. Morgan,²¹⁶ Michael J. Morgan,²¹⁷ Michael J. Morgan,²¹⁸ Michael J. Morgan,²¹⁹ Michael J. Morgan,²²⁰ Michael J. Morgan,²²¹ Michael J. Morgan,²²² Michael J. Morgan,²²³ Michael J. Morgan,²²⁴ Michael J. Morgan,²²⁵ Michael J. Morgan,²²⁶ Michael J. Morgan,²²⁷ Michael J. Morgan,²²⁸ Michael J. Morgan,²²⁹ Michael J. Morgan,²³⁰ Michael J. Morgan,²³¹ Michael J. Morgan,²³² Michael J. Morgan,²³³ Michael J. Morgan,²³⁴ Michael J. Morgan,²³⁵ Michael J. Morgan,²³⁶ Michael J. Morgan,²³⁷ Michael J. Morgan,²³⁸ Michael J. Morgan,²³⁹ Michael J. Morgan,²⁴⁰ Michael J. Morgan,²⁴¹ Michael J. Morgan,²⁴² Michael J. Morgan,²⁴³ Michael J. Morgan,²⁴⁴ Michael J. Morgan,²⁴⁵ Michael J. Morgan,²⁴⁶ Michael J. Morgan,²⁴⁷ Michael J. Morgan,²⁴⁸ Michael J. Morgan,²⁴⁹ Michael J. Morgan,²⁵⁰ Michael J. Morgan,²⁵¹ Michael J. Morgan,²⁵² Michael J. Morgan,²⁵³ Michael J. Morgan,²⁵⁴ Michael J. Morgan,²⁵⁵ Michael J. Morgan,²⁵⁶ Michael J. Morgan,²⁵⁷ Michael J. Morgan,²⁵⁸ Michael J. Morgan,²⁵⁹ Michael J. Morgan,²⁶⁰ Michael J. Morgan,²⁶¹ Michael J. Morgan,²⁶² Michael J. Morgan,²⁶³ Michael J. Morgan,²⁶⁴ Michael J. Morgan,²⁶⁵ Michael J. Morgan,²⁶⁶ Michael J. Morgan,²⁶⁷ Michael J. Morgan,²⁶⁸ Michael J. Morgan,²⁶⁹ Michael J. Morgan,²⁷⁰ Michael J. Morgan,²⁷¹ Michael J. Morgan,²⁷² Michael J. Morgan,²⁷³ Michael J. Morgan,²⁷⁴ Michael J. Morgan,²⁷⁵ Michael J. Morgan,²⁷⁶ Michael J. Morgan,²⁷⁷ Michael J. Morgan,²⁷⁸ Michael J. Morgan,²⁷⁹ Michael J. Morgan,²⁸⁰ Michael J. Morgan,²⁸¹ Michael J. Morgan,²⁸² Michael J. Morgan,²⁸³ Michael J. Morgan,²⁸⁴ Michael J. Morgan,²⁸⁵ Michael J. Morgan,²⁸⁶ Michael J. Morgan,²⁸⁷ Michael J. Morgan,²⁸⁸ Michael J. Morgan,²⁸⁹ Michael J. Morgan,²⁹⁰ Michael J. Morgan,²⁹¹ Michael J. Morgan,²⁹² Michael J. Morgan,²⁹³ Michael J. Morgan,²⁹⁴ Michael J. Morgan,²⁹⁵ Michael J. Morgan,²⁹⁶ Michael J. Morgan,²⁹⁷ Michael J. Morgan,²⁹⁸ Michael J. Morgan,²⁹⁹ Michael J. Morgan,³⁰⁰ Michael J. Morgan,³⁰¹ Michael J. Morgan,³⁰² Michael J. Morgan,³⁰³ Michael J. Morgan,³⁰⁴ Michael J. Morgan,³⁰⁵ Michael J. Morgan,³⁰⁶ Michael J. Morgan,³⁰⁷ Michael J. Morgan,³⁰⁸ Michael J. Morgan,³⁰⁹ Michael J. Morgan,³¹⁰ Michael J. Morgan,³¹¹ Michael J. Morgan,³¹² Michael J. Morgan,³¹³ Michael J. Morgan,³¹⁴ Michael J. Morgan,³¹⁵ Michael J. Morgan,³¹⁶ Michael J. Morgan,³¹⁷ Michael J. Morgan,³¹⁸ Michael J. Morgan,³¹⁹ Michael J. Morgan,³²⁰ Michael J. Morgan,³²¹ Michael J. Morgan,³²² Michael J. Morgan,³²³ Michael J. Morgan,³²⁴ Michael J. Morgan,³²⁵ Michael J. Morgan,³²⁶ Michael J. Morgan,³²⁷ Michael J. Morgan,³²⁸ Michael J. Morgan,³²⁹ Michael J. Morgan,³³⁰ Michael J. Morgan,³³¹ Michael J. Morgan,³³² Michael J. Morgan,³³³ Michael J. Morgan,³³⁴ Michael J. Morgan,³³⁵ Michael J. Morgan,³³⁶ Michael J. Morgan,³³⁷ Michael J. Morgan,³³⁸ Michael J. Morgan,³³⁹ Michael J. Morgan,³⁴⁰ Michael J. Morgan,³⁴¹ Michael J. Morgan,³⁴² Michael J. Morgan,³⁴³ Michael J. Morgan,³⁴⁴ Michael J. Morgan,³⁴⁵ Michael J. Morgan,³⁴⁶ Michael J. Morgan,³⁴⁷ Michael J. Morgan,³⁴⁸ Michael J. Morgan,³⁴⁹ Michael J. Morgan,³⁵⁰ Michael J. Morgan,³⁵¹ Michael J. Morgan,³⁵² Michael J. Morgan,³⁵³ Michael J. Morgan,³⁵⁴ Michael J. Morgan,³⁵⁵ Michael J. Morgan,³⁵⁶ Michael J. Morgan,³⁵⁷ Michael J. Morgan,³⁵⁸ Michael J. Morgan,³⁵⁹ Michael J. Morgan,³⁶⁰ Michael J. Morgan,³⁶¹ Michael J. Morgan,³⁶² Michael J. Morgan,³⁶³ Michael J. Morgan,³⁶⁴ Michael J. Morgan,³⁶⁵ Michael J. Morgan,³⁶⁶ Michael J. Morgan,³⁶⁷ Michael J. Morgan,³⁶⁸ Michael J. Morgan,³⁶⁹ Michael J. Morgan,³⁷⁰ Michael J. Morgan,³⁷¹ Michael J. Morgan,³⁷² Michael J. Morgan,³⁷³ Michael J. Morgan,³⁷⁴ Michael J. Morgan,³⁷⁵ Michael J. Morgan,³⁷⁶ Michael J. Morgan,³⁷⁷ Michael J. Morgan,³⁷⁸ Michael J. Morgan,³⁷⁹ Michael J. Morgan,³⁸⁰ Michael J. Morgan,³⁸¹ Michael J. Morgan,³⁸² Michael J. Morgan,³⁸³ Michael J. Morgan,³⁸⁴ Michael J. Morgan,³⁸⁵ Michael J. Morgan,³⁸⁶ Michael J. Morgan,³⁸⁷ Michael J. Morgan,³⁸⁸ Michael J. Morgan,³⁸⁹ Michael J. Morgan,³⁹⁰ Michael J. Morgan,³⁹¹ Michael J. Morgan,³⁹² Michael J. Morgan,³⁹³ Michael J. Morgan,³⁹⁴ Michael J. Morgan,³⁹⁵ Michael J. Morgan,³⁹⁶ Michael J. Morgan,³⁹⁷ Michael J. Morgan,³⁹⁸ Michael J. Morgan,³⁹⁹ Michael J. Morgan,⁴⁰⁰ Michael J. Morgan,⁴⁰¹ Michael J. Morgan,⁴⁰² Michael J. Morgan,⁴⁰³ Michael J. Morgan,⁴⁰⁴ Michael J. Morgan,⁴⁰⁵ Michael J. Morgan,⁴⁰⁶ Michael J. Morgan,⁴⁰⁷ Michael J. Morgan,⁴⁰⁸ Michael J. Morgan,⁴⁰⁹ Michael J. Morgan,⁴¹⁰ Michael J. Morgan,⁴¹¹ Michael J. Morgan,⁴¹² Michael J. Morgan,⁴¹³ Michael J. Morgan,⁴¹⁴ Michael J. Morgan,⁴¹⁵ Michael J. Morgan,⁴¹⁶ Michael J. Morgan,⁴¹⁷ Michael J. Morgan,⁴¹⁸ Michael J. Morgan,⁴¹⁹ Michael J. Morgan,⁴²⁰ Michael J. Morgan,⁴²¹ Michael J. Morgan,⁴²² Michael J. Morgan,⁴²³ Michael J. Morgan,⁴²⁴ Michael J. Morgan,⁴²⁵ Michael J. Morgan,⁴²⁶ Michael J. Morgan,⁴²⁷ Michael J. Morgan,⁴²⁸ Michael J. Morgan,⁴²⁹ Michael J. Morgan,⁴³⁰ Michael J. Morgan,⁴³¹ Michael J. Morgan,⁴³² Michael J. Morgan,⁴³³ Michael J. Morgan,⁴³⁴ Michael J. Morgan,⁴³⁵ Michael J. Morgan,⁴³⁶ Michael J. Morgan,⁴³⁷ Michael J. Morgan,⁴³⁸ Michael J. Morgan,⁴³⁹ Michael J. Morgan,⁴⁴⁰ Michael J. Morgan,⁴⁴¹ Michael J. Morgan,⁴⁴² Michael J. Morgan,⁴⁴³ Michael J. Morgan,⁴⁴⁴ Michael J. Morgan,⁴⁴⁵ Michael J. Morgan,⁴⁴⁶ Michael J. Morgan,⁴⁴⁷ Michael J. Morgan,⁴⁴⁸ Michael J. Morgan,⁴⁴⁹ Michael J. Morgan,⁴⁵⁰ Michael J. Morgan,⁴⁵¹ Michael J. Morgan,⁴⁵² Michael J. Morgan,⁴⁵³ Michael J. Morgan,⁴⁵⁴ Michael J. Morgan,⁴⁵⁵ Michael J. Morgan,⁴⁵⁶ Michael J. Morgan,⁴⁵⁷ Michael J. Morgan,⁴⁵⁸ Michael J. Morgan,⁴⁵⁹ Michael J. Morgan,⁴⁶⁰ Michael J. Morgan,⁴⁶¹ Michael J. Morgan,⁴⁶² Michael J. Morgan,⁴⁶³ Michael J. Morgan,⁴⁶⁴ Michael J. Morgan,⁴⁶⁵ Michael J. Morgan,⁴⁶⁶ Michael J. Morgan,⁴⁶⁷ Michael J. Morgan,⁴⁶⁸ Michael J. Morgan,⁴⁶⁹ Michael J. Morgan,⁴⁷⁰ Michael J. Morgan,⁴⁷¹ Michael J. Morgan,⁴⁷² Michael J. Morgan,⁴⁷³ Michael J. Morgan,⁴⁷⁴ Michael J. Morgan,⁴⁷⁵ Michael J. Morgan,⁴⁷⁶ Michael J. Morgan,⁴⁷⁷ Michael J. Morgan,⁴⁷⁸ Michael J. Morgan,⁴⁷⁹ Michael J. Morgan,⁴⁸⁰ Michael J. Morgan,⁴⁸¹ Michael J. Morgan,⁴⁸² Michael J. Morgan,⁴⁸³ Michael J. Morgan,⁴⁸⁴ Michael J. Morgan,⁴⁸⁵ Michael J. Morgan,⁴⁸⁶ Michael J. Morgan,⁴⁸⁷ Michael J. Morgan,⁴⁸⁸ Michael J. Morgan,⁴⁸⁹ Michael J. Morgan,⁴⁹⁰ Michael J. Morgan,⁴⁹¹ Michael J. Morgan,⁴⁹² Michael J. Morgan,⁴⁹³ Michael J. Morgan,⁴⁹⁴ Michael J. Morgan,⁴⁹⁵ Michael J. Morgan,⁴⁹⁶ Michael J. Morgan,⁴⁹⁷ Michael J. Morgan,⁴⁹⁸ Michael J. Morgan,⁴⁹⁹ Michael J. Morgan,⁵⁰⁰ Michael J. Morgan,⁵⁰¹ Michael J. Morgan,⁵⁰² Michael J. Morgan,⁵⁰³ Michael J. Morgan,⁵⁰⁴ Michael J. Morgan,⁵⁰⁵ Michael J. Morgan,⁵⁰⁶ Michael J. Morgan,⁵⁰⁷ Michael J. Morgan,⁵⁰⁸ Michael J. Morgan,⁵⁰⁹ Michael J. Morgan,⁵¹⁰ Michael J. Morgan,⁵¹¹ Michael J. Morgan,⁵¹² Michael J. Morgan,⁵¹³ Michael J. Morgan,⁵¹⁴ Michael J. Morgan,⁵¹⁵ Michael J. Morgan,⁵¹⁶ Michael J. Morgan,⁵¹⁷ Michael J. Morgan,⁵¹⁸ Michael J. Morgan,⁵¹⁹ Michael J. Morgan,⁵²⁰ Michael J. Morgan,⁵²¹ Michael J. Morgan,⁵²² Michael J. Morgan,⁵²³ Michael J. Morgan,⁵²⁴ Michael J. Morgan,⁵²⁵ Michael J. Morgan,⁵²⁶ Michael J. Morgan,⁵²⁷ Michael J. Morgan,⁵²⁸ Michael J. Morgan,⁵²⁹ Michael J. Morgan,⁵³⁰ Michael J. Morgan,⁵³¹ Michael J. Morgan,⁵³² Michael J. Morgan,⁵³³ Michael J. Morgan,⁵³⁴ Michael J. Morgan,⁵³⁵ Michael J. Morgan,⁵³⁶ Michael J. Morgan,⁵³⁷ Michael J. Morgan,⁵³⁸ Michael J. Morgan,⁵³⁹ Michael J. Morgan,⁵⁴⁰ Michael J. Morgan,⁵⁴¹ Michael J. Morgan,⁵⁴² Michael J. Morgan,⁵⁴³ Michael J. Morgan,⁵⁴⁴ Michael J. Morgan,⁵⁴⁵ Michael J. Morgan,⁵⁴⁶ Michael J. Morgan,⁵⁴⁷ Michael J. Morgan,⁵⁴⁸ Michael J. Morgan,⁵⁴⁹ Michael J. Morgan,⁵⁵⁰ Michael J. Morgan,⁵⁵¹ Michael J. Morgan,⁵⁵² Michael J. Morgan,⁵⁵³ Michael J. Morgan,⁵⁵⁴ Michael J. Morgan,⁵⁵⁵ Michael J. Morgan,⁵⁵⁶ Michael J. Morgan,⁵⁵⁷ Michael J. Morgan,⁵⁵⁸ Michael J. Morgan,⁵⁵⁹ Michael J. Morgan,⁵⁶⁰ Michael J. Morgan,⁵⁶¹ Michael J. Morgan,⁵⁶² Michael J. Morgan,⁵⁶³ Michael J. Morgan,⁵⁶⁴ Michael J. Morgan,⁵⁶⁵ Michael J. Morgan,⁵⁶⁶ Michael J. Morgan,⁵⁶⁷ Michael J. Morgan,⁵⁶⁸ Michael J. Morgan,⁵⁶⁹ Michael J. Morgan,⁵⁷⁰ Michael J. Morgan,⁵⁷¹ Michael J. Morgan,⁵⁷² Michael J. Morgan,⁵⁷³ Michael J. Morgan,⁵⁷⁴ Michael J. Morgan,⁵⁷⁵ Michael J. Morgan,⁵⁷⁶ Michael J. Morgan,⁵⁷⁷ Michael J. Morgan,⁵⁷⁸ Michael J. Morgan,⁵⁷⁹ Michael J. Morgan,⁵⁸⁰ Michael J. Morgan,⁵⁸¹ Michael J. Morgan,⁵⁸² Michael J. Morgan,⁵⁸³ Michael J. Morgan,⁵⁸⁴ Michael J. Morgan,⁵⁸⁵ Michael J. Morgan,⁵⁸⁶ Michael J. Morgan,⁵⁸⁷ Michael J. Morgan,⁵⁸⁸ Michael J. Morgan,⁵⁸⁹ Michael J. Morgan,⁵⁹⁰ Michael J. Morgan,⁵⁹¹ Michael J. Morgan,⁵⁹² Michael J. Morgan,⁵⁹³ Michael J. Morgan,⁵⁹⁴ Michael J. Morgan,⁵⁹⁵ Michael J. Morgan,⁵⁹⁶ Michael J. Morgan,⁵⁹⁷ Michael J. Morgan,⁵⁹⁸ Michael J. Morgan,⁵⁹⁹ Michael J. Morgan,⁶⁰⁰ Michael J. Morgan,⁶⁰¹ Michael J. Morgan,⁶⁰² Michael J. Morgan,⁶⁰³ Michael J. Morgan,⁶⁰⁴ Michael J. Morgan,⁶⁰⁵ Michael J. Morgan,⁶⁰⁶ Michael J. Morgan,⁶⁰⁷ Michael J. Morgan,⁶⁰⁸ Michael J. Morgan,⁶⁰⁹ Michael J. Morgan,⁶¹⁰ Michael J. Morgan,⁶¹¹ Michael J. Morgan,⁶¹² Michael J. Morgan,⁶¹³ Michael J. Morgan,⁶¹⁴ Michael J. Morgan,⁶¹⁵ Michael J. Morgan,⁶¹⁶ Michael J. Morgan,⁶¹⁷ Michael J. Morgan,⁶¹⁸ Michael J. Morgan,⁶¹⁹ Michael J. Morgan,⁶²⁰ Michael J. Morgan,⁶²¹ Michael J. Morgan,⁶²² Michael J. Morgan,⁶²³ Michael J. Morgan,⁶²⁴ Michael J. Morgan,⁶²⁵ Michael J. Morgan,⁶²⁶ Michael J. Morgan,⁶²⁷ Michael J. Morgan,⁶²⁸ Michael J. Morgan,⁶²⁹ Michael J. Morgan,⁶³⁰ Michael J. Morgan,⁶³¹ Michael J. Morgan,⁶³² Michael J. Morgan,⁶³³ Michael J. Morgan,⁶³⁴ Michael J. Morgan,⁶³⁵ Michael J. Morgan,⁶³⁶ Michael J. Morgan,⁶³⁷ Michael J. Morgan,⁶³⁸ Michael J. Morgan,⁶³⁹ Michael J. Morgan,⁶⁴⁰ Michael J. Morgan,⁶⁴¹ Michael J. Morgan,⁶⁴² Michael J. Morgan,⁶⁴³ Michael J. Morgan,⁶⁴⁴ Michael J. Morgan,⁶⁴⁵ Michael J. Morgan,⁶⁴⁶ Michael J. Morgan,⁶⁴⁷ Michael J. Morgan,⁶⁴⁸ Michael J. Morgan,⁶⁴⁹ Michael J. Morgan,⁶⁵⁰ Michael J. Morgan,⁶⁵¹ Michael J. Morgan,⁶⁵² Michael J. Morgan,⁶⁵³ Michael J. Morgan,⁶⁵⁴ Michael J. Morgan,⁶⁵⁵ Michael J. Morgan,⁶⁵⁶ Michael J. Morgan,⁶⁵⁷ Michael J. Morgan,⁶⁵⁸ Michael J. Morgan,⁶⁵⁹ Michael J. Morgan,⁶⁶⁰ Michael J. Morgan,⁶⁶¹ Michael J. Morgan,⁶⁶² Michael J. Morgan,⁶⁶³ Michael J. Morgan,⁶⁶⁴ Michael J. Morgan,⁶⁶⁵ Michael J. Morgan,⁶⁶⁶ Michael J. Morgan,⁶⁶⁷ Michael J. Morgan,⁶⁶⁸ Michael J. Morgan,⁶⁶⁹ Michael J. Morgan,⁶⁷⁰ Michael J. Morgan,⁶⁷¹ Michael J. Morgan,⁶⁷² Michael J. Morgan,⁶⁷³ Michael J. Morgan,⁶⁷⁴ Michael J. Morgan,⁶⁷⁵ Michael J. Morgan,⁶⁷⁶ Michael J. Morgan,⁶⁷⁷ Michael J. Morgan,⁶⁷⁸ Michael J. Morgan,⁶⁷⁹ Michael J. Morgan,⁶⁸⁰ Michael J. Morgan,⁶⁸¹ Michael J. Morgan,⁶⁸² Michael J. Morgan,⁶⁸³ Michael J. Morgan,⁶⁸⁴ Michael J. Morgan,⁶⁸⁵ Michael J. Morgan,⁶⁸⁶ Michael J. Morgan,⁶⁸⁷ Michael J. Morgan,⁶⁸⁸ Michael J. Morgan,⁶⁸⁹ Michael J. Morgan,⁶⁹⁰ Michael J. Morgan,⁶⁹¹ Michael J. Morgan,⁶⁹² Michael J. Morgan,⁶⁹³ Michael J. Morgan,⁶⁹⁴ Michael J. Morgan,⁶⁹⁵ Michael J. Morgan,⁶⁹⁶ Michael J. Morgan,⁶⁹⁷ Michael J. Morgan,⁶⁹⁸ Michael J. Morgan,⁶⁹⁹ Michael J. Morgan,⁷⁰⁰ Michael J. Morgan,⁷⁰¹ Michael J. Morgan,⁷⁰² Michael J. Morgan,⁷⁰³ Michael J. Morgan,⁷⁰⁴ Michael J. Morgan,⁷⁰⁵ Michael J. Morgan,⁷⁰⁶ Michael J. Morgan,⁷⁰⁷ Michael J. Morgan,⁷⁰⁸ Michael J. Morgan,⁷⁰⁹ Michael J. Morgan,⁷¹⁰ Michael J. Morgan,⁷¹¹ Michael J. Morgan,⁷¹² Michael J. Morgan,⁷¹³ Michael J. Morgan,⁷¹⁴ Michael J. Morgan,⁷¹⁵ Michael J. Morgan,⁷¹⁶ Michael J. Morgan,⁷¹⁷ Michael J. Morgan,⁷¹⁸ Michael J. Morgan,⁷¹⁹ Michael J. Morgan,⁷²⁰ Michael J. Morgan,⁷²¹ Michael J. Morgan,⁷²² Michael J. Morgan,⁷²³ Michael J. Morgan,⁷²⁴ Michael J. Morgan,⁷²⁵ Michael J. Morgan,⁷²⁶ Michael J. Morgan,⁷²⁷ Michael J. Morgan,⁷²⁸ Michael J. Morgan,⁷²⁹ Michael J. Morgan,⁷³⁰ Michael J. Morgan,⁷³¹ Michael J. Morgan,⁷³² Michael J. Morgan,⁷³³ Michael J. Morgan,⁷³⁴ Michael J. Morgan,⁷³⁵ Michael J. Morgan,⁷³⁶ Michael J. Morgan,⁷³⁷ Michael J. Morgan,⁷³⁸ Michael J. Morgan,⁷³⁹ Michael J. Morgan,⁷⁴⁰ Michael J. Morgan,⁷⁴¹ Michael J. Morgan,⁷⁴² Michael J. Morgan,⁷⁴³ Michael J. Morgan,⁷⁴⁴ Michael J. Morgan,⁷⁴⁵ Michael J. Morgan,⁷⁴⁶ Michael J. Morgan,⁷⁴⁷ Michael J. Morgan,⁷⁴⁸ Michael J. Morgan,⁷⁴⁹ Michael J. Morgan,⁷⁵⁰ Michael J. Morgan,⁷⁵¹ Michael J. Morgan,⁷⁵² Michael J. Morgan,⁷⁵³ Michael J. Morgan,⁷⁵⁴ Michael J. Morgan,⁷⁵⁵ Michael J. Morgan,⁷⁵⁶ Michael J. Morgan,⁷⁵⁷ Michael J. Morgan,⁷⁵⁸ Michael J. Morgan,⁷⁵⁹ Michael J. Morgan,⁷⁶⁰ Michael J. Morgan,⁷⁶¹ Michael J. Morgan,⁷⁶² Michael J. Morgan,⁷⁶³ Michael J. Morgan,⁷⁶⁴ Michael J. Morgan,⁷⁶⁵ Michael J. Morgan,⁷⁶⁶ Michael J. Morgan,⁷⁶⁷ Michael J. Morgan,⁷⁶⁸ Michael J. Morgan



nature

Vol 437|15 September 2005|doi:10.1038/nature03959

ARTICLES

Genome sequencing in microfabricated high-density picolitre reactors

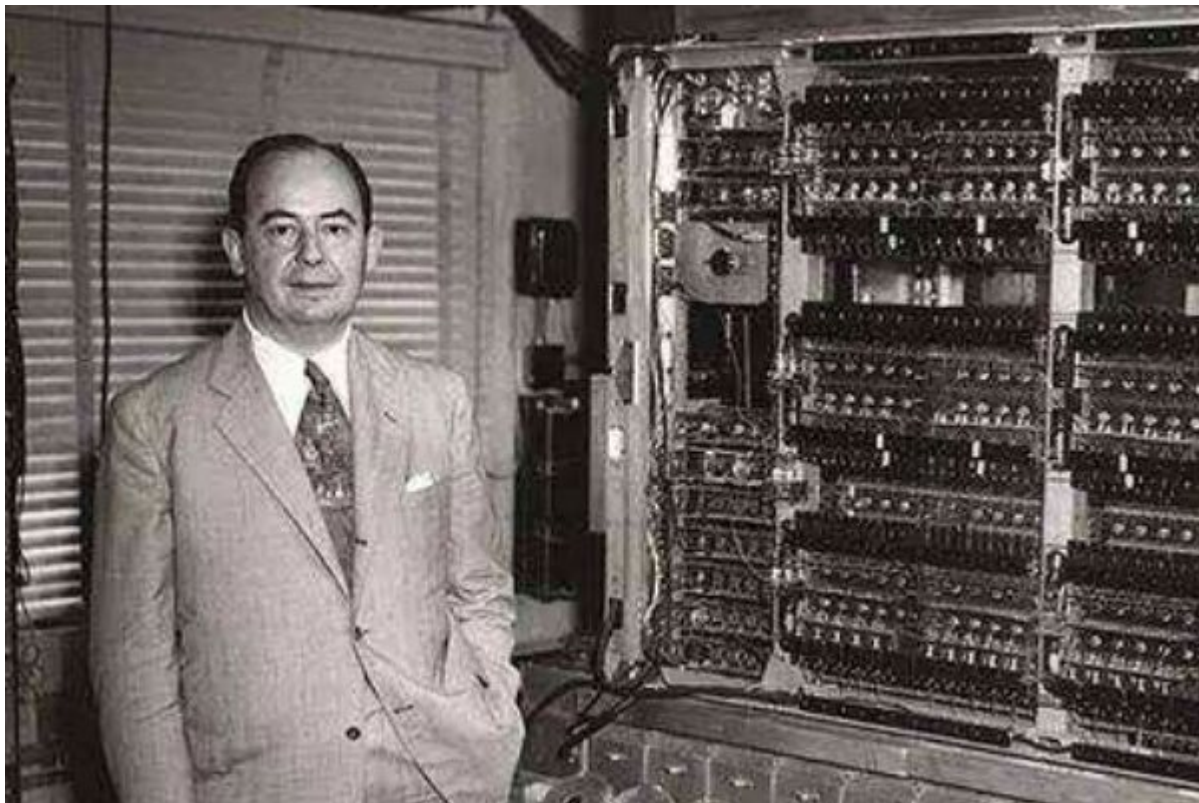
Marcel Margulies^{1*}, Michael Egholm^{1*}, William E. Altman¹, Said Attiya¹, Joel S. Bader¹, Lisa A. Bemben¹, Jan Berka¹, Michael S. Braverman¹, Yi-Ju Chen¹, Zhoutao Chen¹, Scott B. Dewell¹, Lei Du¹, Joseph M. Fierro¹, Xavier V. Gomez¹, David H. Hays¹, John E. Healy¹, Steven M. Leake¹, Ming Li¹, Yong Li¹, Kenton L. Lohman¹, Hong Lu¹, Vinod B. Makhijani¹, Keith E. McDade¹, Michael P. McKenna¹, Eugene W. Myers², Elizabeth Nickerson¹, John R. Nobile¹, Ramona Plant¹, Bernard P. Puc¹, Michael T. Ronan¹, George T. Roth¹, Gary J. Sarkis¹, Jan Fredrik Simons¹, John W. Simpson¹, Maithreyan Srinivasan¹, Karrie R. Tartaro¹, Alexander Tomasz³, Kari A. Vogt¹, Greg A. Volkmer¹, Shally H. Wang¹, Yong Wang¹, Michael P. Weiner⁴, Pengguang Yu¹, Richard F. Begley¹ & Jonathan M. Rothberg¹

2005年，新一代测序技术出现。

The proliferation of large-scale DNA-sequencing projects in recent years has driven a search for alternative methods to reduce time and cost. Here we describe a scalable, highly parallel sequencing system with raw throughput significantly greater than that of state-of-the-art capillary electrophoresis instruments. The apparatus uses a novel fibre-optic slide of individual wells and is able to sequence 25 million bases, at 99% or better accuracy, in one four-hour run. To achieve an approximately 100-fold increase in throughput over current Sanger sequencing technology, we have developed an emulsion method for DNA amplification and an instrument for sequencing by synthesis using a pyrosequencing protocol



约翰·冯诺依曼和计算机

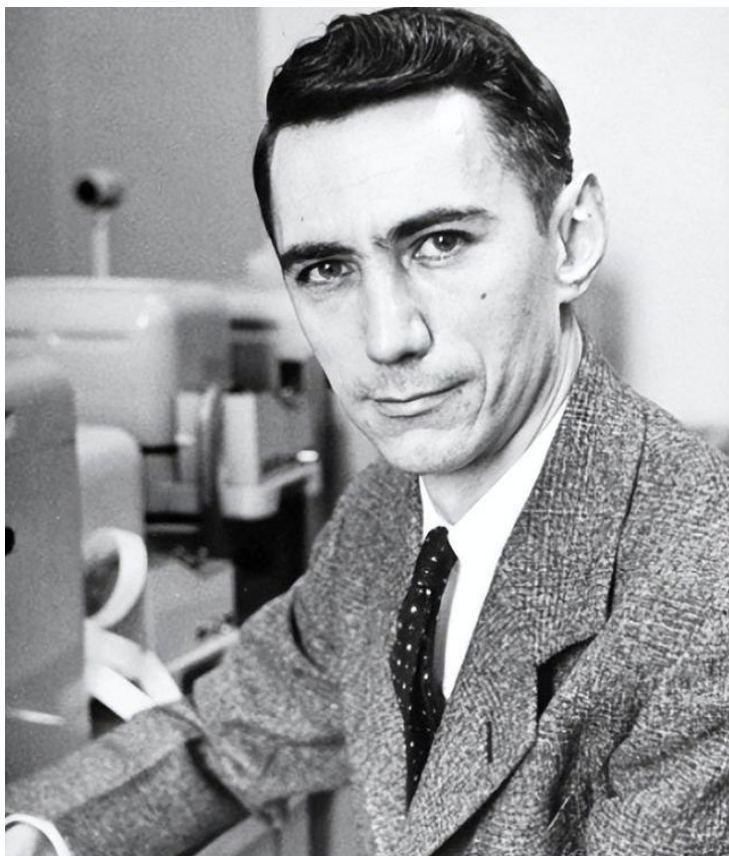


- 算符环理论
- 博弈论
- 蒙特卡洛方法
- 冯诺依曼体系

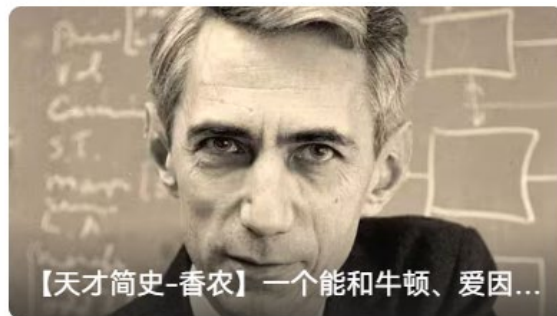
计算机制造的三个基本原则，即采用二进制逻辑、程序存储执行以及计算机由五个部分组成(运算器、控制器、存储器、输入设备、输出设备)



信息论之父: 克劳德·艾尔伍德·香农



$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log p(x)$$



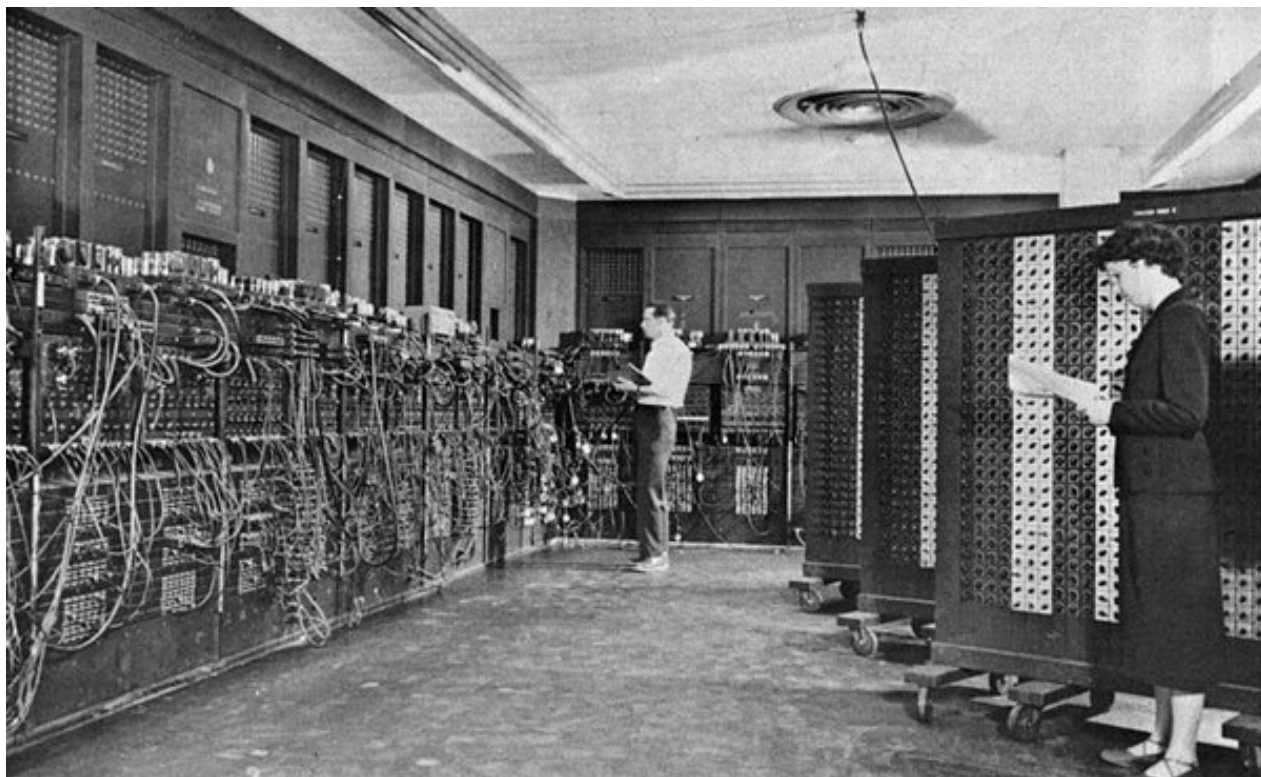


计算机科学与人工智能之父: 艾伦·麦席森·图灵





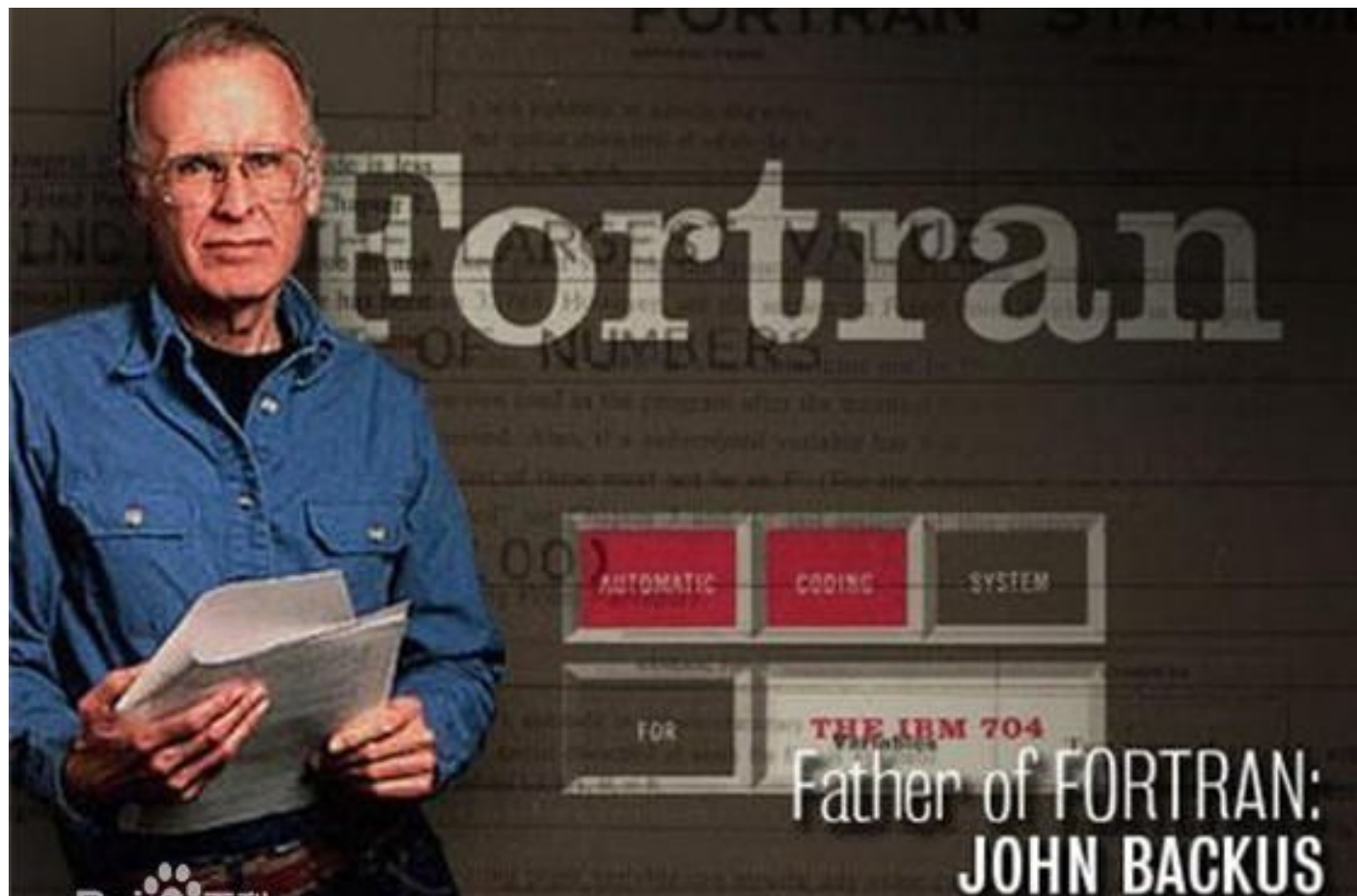
1946年，第一台计算机，ENIAC（埃尼亚克）



ENIAC：长30.48米，宽1米，占地面积约63平方米，30个操作台，约相当于10间普通房间的大小，重达30吨，耗电量150千瓦，造价48万美元。它包含了17,468 真空管7,200水晶 二极管, 1,500 中转, 70,000 电阻器, 10,000 电容器, 1500 继电器, 6000多个开关，每秒执行5000次加法或400次乘法，是继电器计算机的1000倍、手工计算的20万倍。



1951年，IBM公司的约翰·贝克斯在纽约正式对外发布Fortran语言，称为FORTRAN I。





1969-1991: WWW的诞生

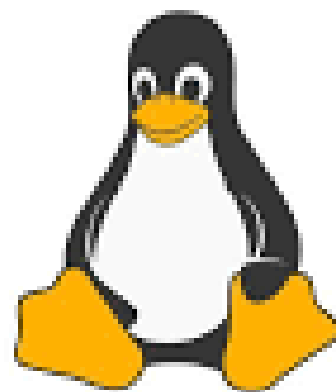
- 1969年，ARPANet建立，对计算机网络技术发展做出重要贡献。
- 1970年后，出现了E-mail、Ethernet、TCP协议。
- 1980年后，以IBM为代表的个人计算机开始普及。
- 1991年，World Wide Web协议被建立。



1991: Linux操作系统



Linus Torvalds



Linux

Shell
Scripting





1952年，肯德鲁用计算机程序来解析蛋白结构。（acta Cryst, 1952）

1962年，**Dayhoff**开发序列分析软件COMPROTEIN。

1965年，**Dayhoff**出版一个蛋白质数据库Atlas（第一年65条序列），发展为1983年的PIR。

1966年，**Dayhoff**对蛋白质家族进化深入研究。（ Science, 1966 ）

1967年，Fitch发表系统发育树（ Science, 1967 ）

1970年，Hesper提出“Bioinformatics”单词，生物信息学概念被定义。

1970年，Needleman和Wunsch提出全局比对算法。（J. Mol. Biol., 1970）

1977年，Protein Data Bank (PDB) 数据库建立。

1978年，**Dayhoff**提出氨基酸序列比对的PAM矩阵。



Dr. Margaret Oakley Dayhoff The Mother of Bioinformatics





1981年，Smith和Waterman发表**局部比对算法**。（J. Mol. Biol., 1981）

1982年，建立核酸序列数据库**Genbank**（最初606条序列）。

1983年，建立Protein information Resource (PIR) 蛋白数据库。2003年（4200万）

1988年，Lipman和Pearson发表FastA算法。（PNAS, 1988）

1990年，Altschul发表**Blast算法**。（J. Mol. Biol., 1990）

1997年，Altschul发表Gapped BLAST和PSI-BLAST算法。（Nucleic Acids Research）

1997年，Chris Burge等发明了GENSCAN算法。（J. Mol. Biol., 1997）

2002年，Kent建立BLAT算法。（Genome Research, 2002）

2003年，Ouzounis对前期的生物信息学发展进行了总结。（BIOINFORMATICS, 2003）



生物信息学发展的四个阶段

- (1) 萌芽期(60-70年代): 以Dayhoff的替换矩阵和Neelleman-Wunsch算法为代表, 它们实际组成了生物信息学的一个最基本的内容和思路: 序列比较。它们的出现, 代表了生物信息学的诞生(虽然“生物信息学”一词很晚才出现);
- (2) 形成期(80年代): 以分子数据库和FASTA等相似性搜索程序为代表。在这一阶段, 生物信息学作为一个新兴学科已经形成, 并确立了自身学科的特征和地位;
- (3) 基因组测序时代(90年代-至今): 以模式基因组测序与BLAST为代表;
- (4) 高通量测序时代(2005 -至今): 以第二和三代测序技术和基因组重测序为代表。



1994年，中国终于获准加入互联网，并在同年5月完成全部中国联网工作。

1998年，中国人类基因组学南方中心和北方中心分别在上海和北京成立。

1999年，华大基因在北京成立。

2003年，中科院北京基因组所成立。

2008年，魏丽萍老师发表了中国生物信息学发展情况。（PLoS Computational Biology, 2008）



我国生物信息学的开拓者



郝柏林 院士
理论物理所 进化发育分析



陈润生 院士
生物物理所 ncRNA



张春霆 院士
天津大学 Z曲线DNA分析



李衍达 院士
清华大学 基因表达调控



孙之荣 教授
清华大学 分子网络分析



罗辽复 教授
内蒙古大学 基因组进化



山東農業大學

SHANDONG AGRICULTURAL UNIVERSITY

生物信息学

第一章: 导论

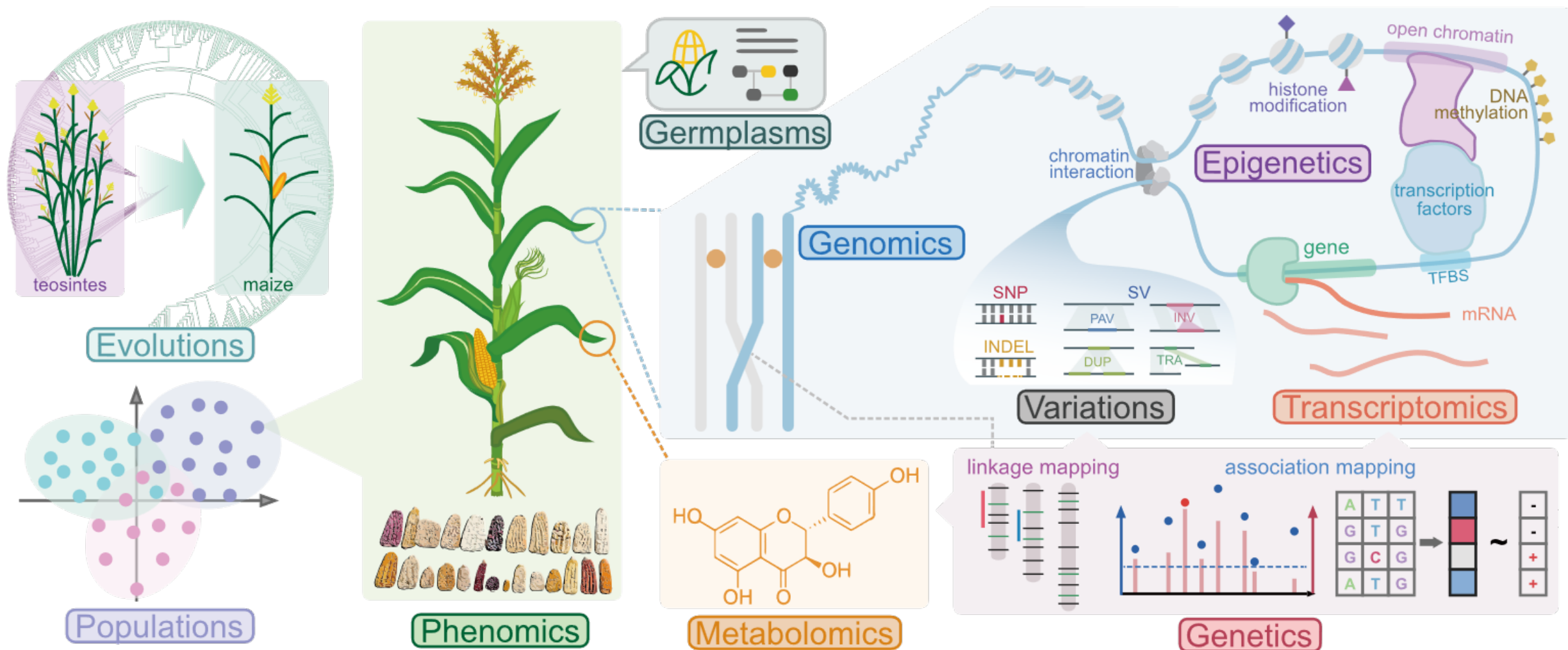
3. 生物信息学研究内容

桂松涛

songtaogui@163.com

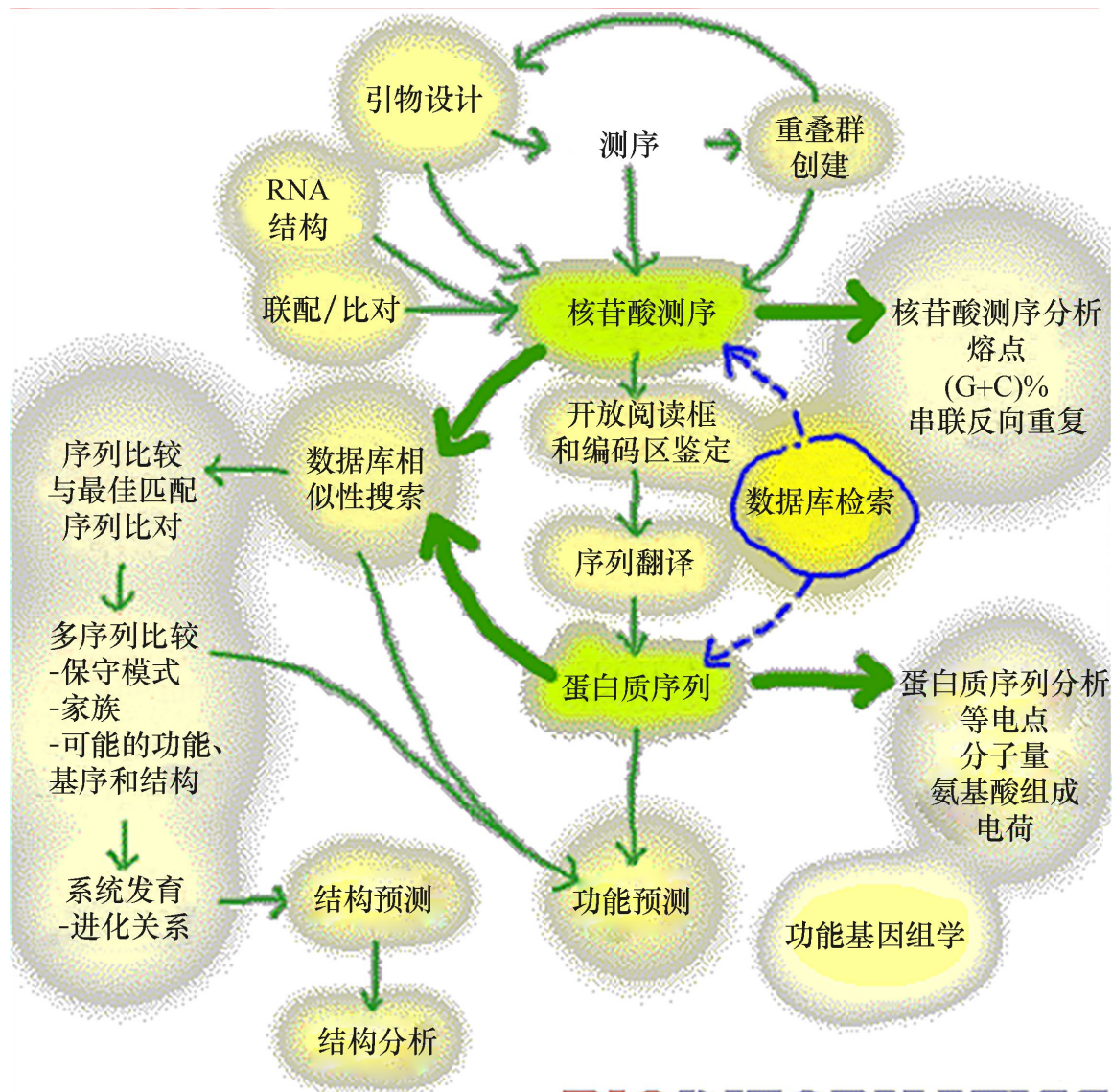


Overview





生物信息学早期研究路线图





Bioinformatics \approx Biology + Programming + Mathematics

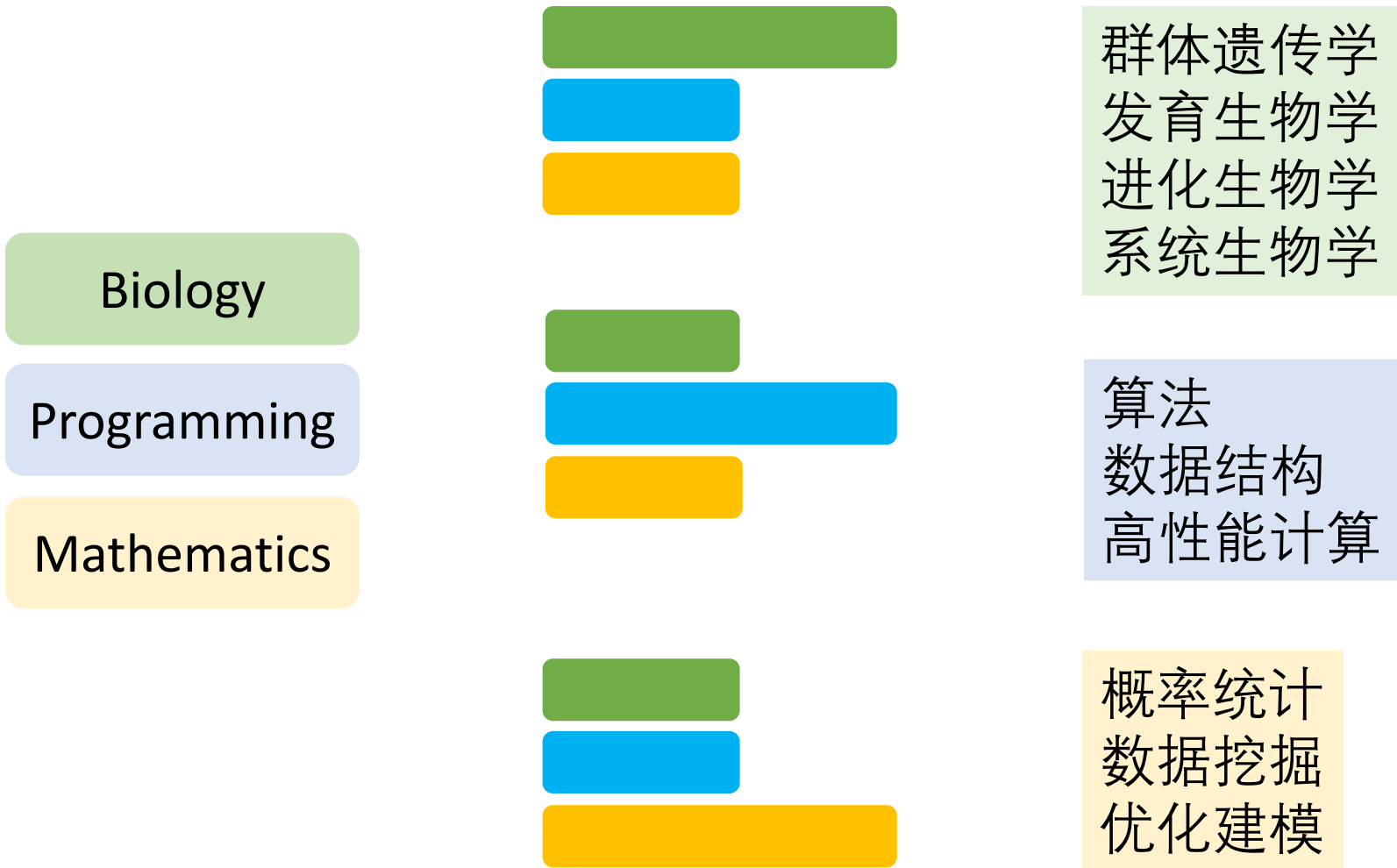
新手村技能加点

1. 基本的 **生物学和组学** 知识: 测序(二代/三代), RNA-Seq, 甲基化, ChIP-seq, ATAC-seq, 单细胞 ...
2. 了解常用的 **生物数据库** 资源: NCBI, Ensemble, UCSC, Uniprot, CoGe, Panther ...
3. **Linux系统** 相关: Shell(Bash)编程, 作业调度, 超算使用 ...
4. 掌握至少一门 **通用编程语言** : python, perl, R, julia, rust, java, C++ ...
5. 基本的 **数据分析** 技巧: 表格操作, 数据可视化 ...
6. 基本的 **统计分析** 知识: 数据分布、假设检验、线性回归、聚类
7. 常见的 **生物信息学算法** : 动态规划, 遗传算法, BWT转换, 序列比对 ...

延伸阅读: **生物信息等级划分**



生物信息进阶转职路线





生信工具使用者的加点建议

我以湿实验为主,偶尔使用常规分析流程、现成的生信软件做一些基础分析,我要学啥?

- 基本的linux操作;
- 对组学技术和测序建库原理有本的认识;
- 了解常用的生信数据库和软件;
- 善用搜索工具找到自己需要的软件,并能看懂说明书;
- 学会向生信专业人员描述问题,快速获得帮助;

"切,电脑里又种不出粮食"

- 编程可以快速训练抽象的能力;
- 生信从业者情绪不稳定(久坐+熬夜+头发少),对我们友好一些 😊;



- 快速正反馈
- 终身学习
- 头脑风暴
- 学好数理化
- 科学问题驱动



程序员有三种美德: 懒惰, 急躁和傲慢

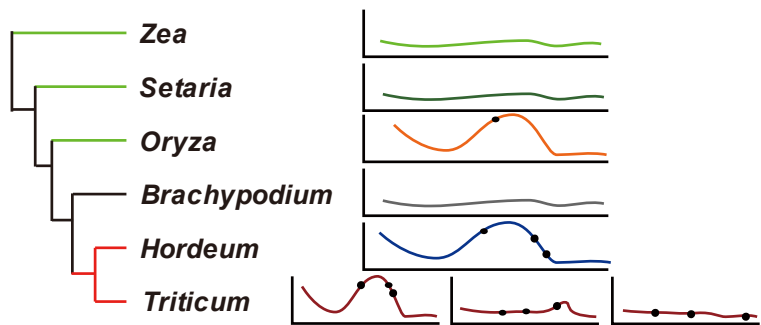
---- Larry Wall



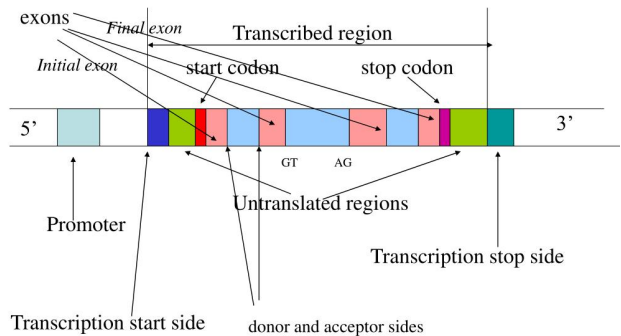
生物信息数据库及其应用

		A	T	G	C	T
	0	-2	-4	-6	-8	-10
A	-2	1	-1	-3	-5	-7
G	-4	-1	0	0	-2	-4
C	-6	-3	-2	-1	1	-1
T	-8	-5	-2	-3	-1	2

序列比对算法



进化分析方法



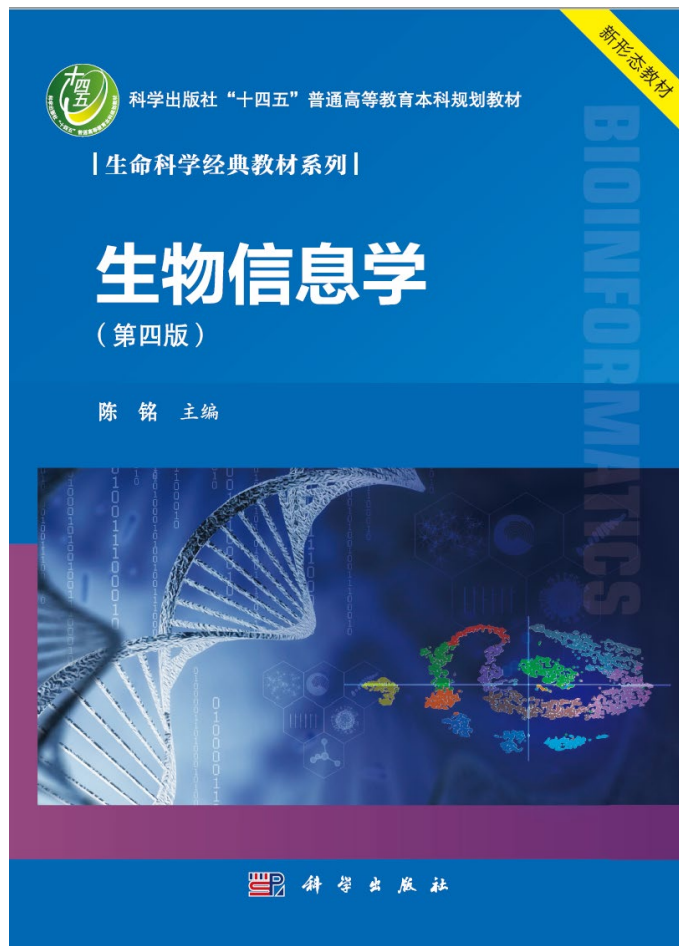
核酸序列分析



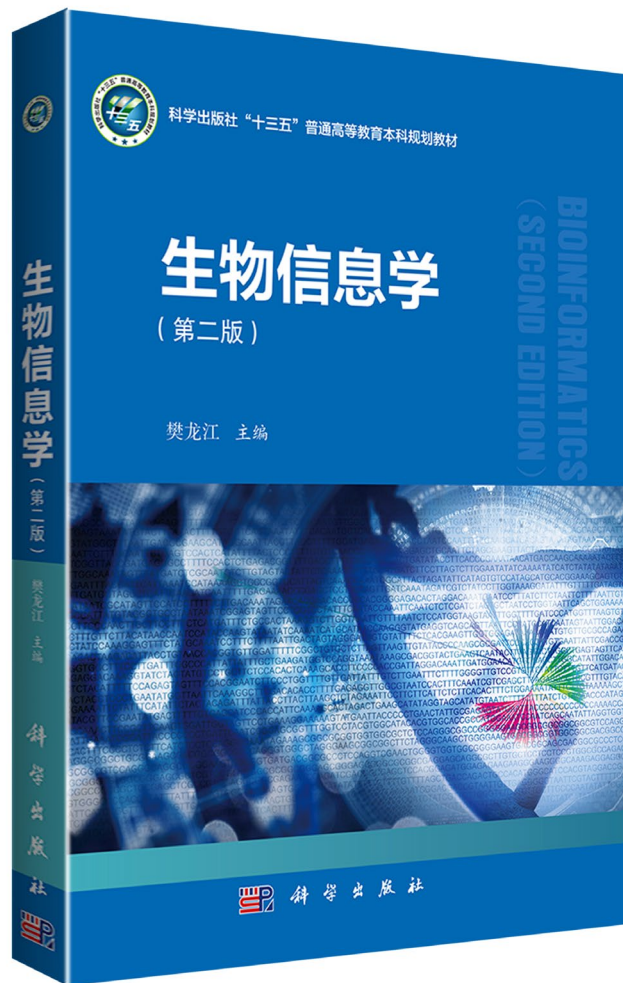
理论课部分（16课时）

1. 生物信息学导论（2课时）
2. 生物信息数据库（2课时）
3. 序列比对（4课时）
4. 进化分析（4课时）
5. 核酸序列分析（4课时）

实验课部分（16课时）



陈铭 主编
科学出版社, 2022



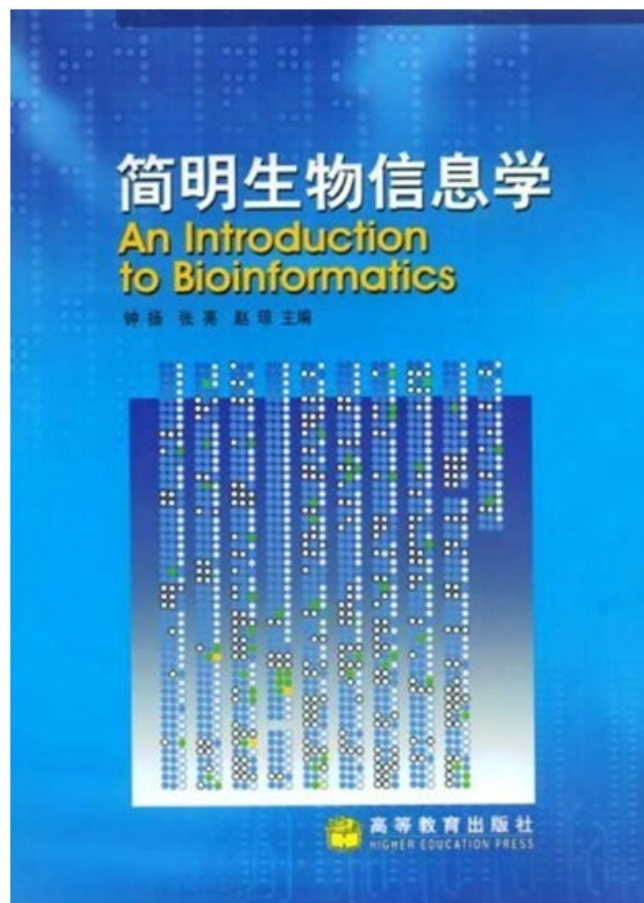
樊龙江 主编
科学出版社, 2021



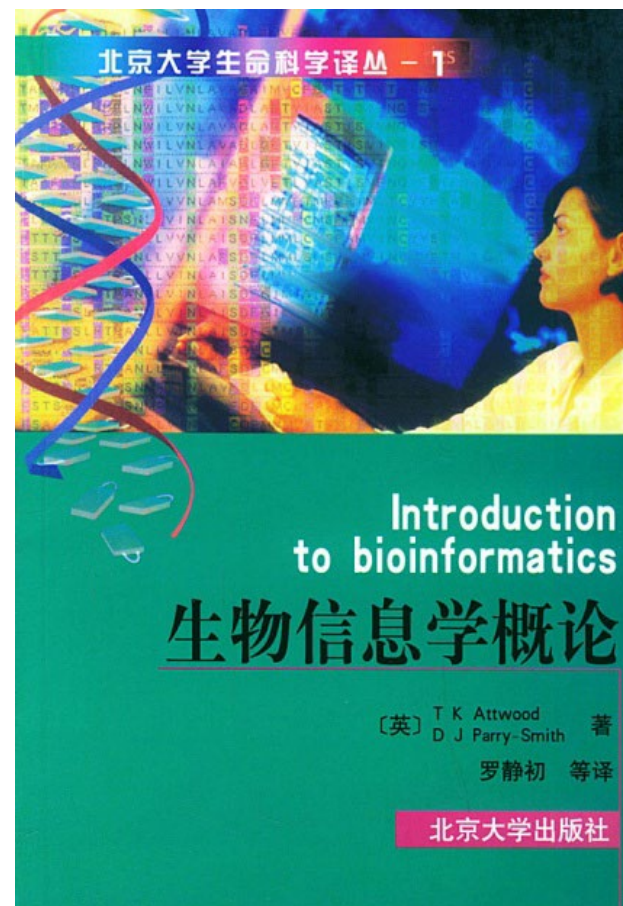
吕巍 主编
高等教育出版社



参考书推荐



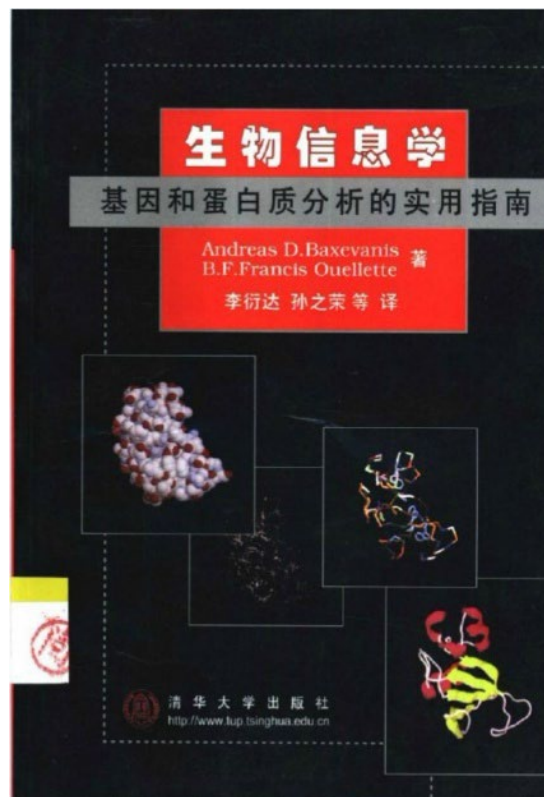
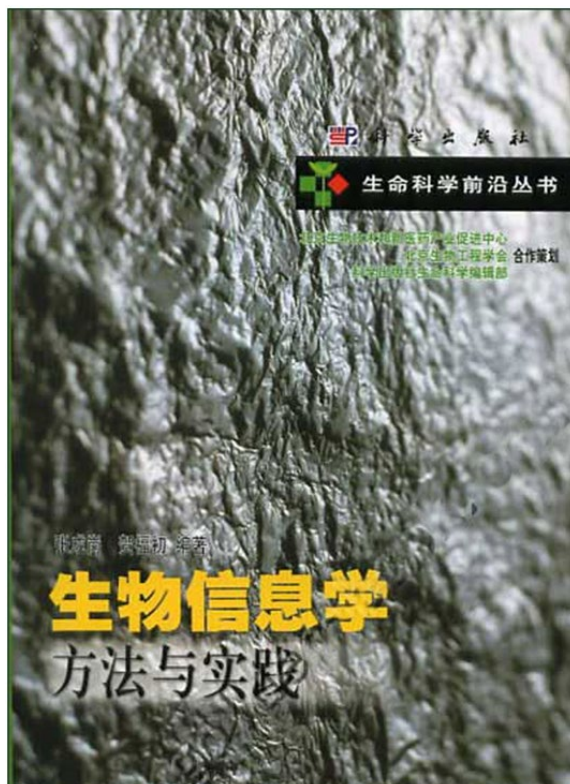
《简明生物信息学》
钟扬，张亮，赵琼 主编
高等教育出版社



《生物信息学概论》
罗静初 等译
北京大学出版社



参考书推荐





平时表现(出勤+测验) : 20%

期末测试: 80%

绝知须躬行

- 课上内容以框架、大纲为主;
- 提供自学资源、资料;
- 90%以上的内容需要课下**自学**掌握;



山東農業大學

SHANDONG AGRICULTURAL UNIVERSITY

生物信息学

Thanks!

桂松涛

山东农业大学 生命科学学院

songtaogui@163.com