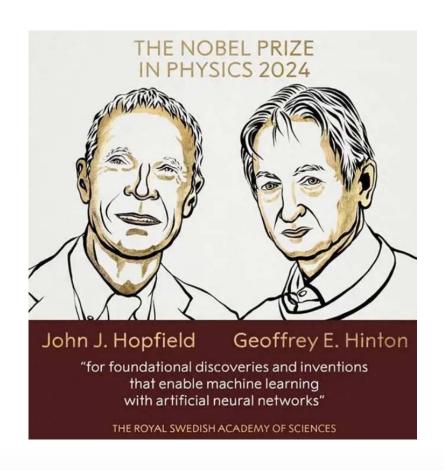


# 生物信息学

### 第一章 历史简介和概论

### 2024年诺贝尔物理学奖 & 化学奖







"基于人工神经网络实现机器学

习的基础性发现和发明"

"计算蛋白质设计"、"蛋白质结 构预测"

Bioinformatics, 2025, HUST

### 概率模型 *vs.* 语言模型



Primer

https://doi.org/10.1038/s41587-024-02123-4

#### Designing proteins with language models

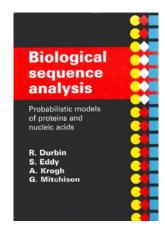
#### Foundations of protein language models

Fundamentally, protein language models aim to predict how likely we are to observe a particular protein sequence S given all the protein sequence data collected thus far. We denote a protein sequence  $S = (s_1, s_2, \dots s_N)$ , where  $s_i$  represents the amino acid at position i in the sequence. As a first approximation, we might consider the probability of observing a protein as the joint probability of observing each of its constituent amino acids. Under this model, referred to as unigram, we calculate the probability of a sequence S as

$$P(S) = \prod_{i}^{N} P(s_i)$$

In practice, to compute P(S), we simply tabulate the frequency of each amino acid occurring in our sequence database and multiply the probabilities for the specific sequence S. However, proteins are not unordered collections of amino acids, Rather, the specific order in which we observe the amino acids is a critical determinant of structure and function. To capture this order dependency, we can use the preceding residues to inform the probability of the next amino acid. In an n-gram model, we multiply these contextualized probabilities to form the overall probability of the sequence:

$$P(S) = \prod_{i=1}^{N} P(s_{i}|s_{i-(n-1)}, \dots, s_{i-1})$$



生物信息学 1.0

AI生物学/生物信息学2.0

Nat Biotechnol, 2024, 42, 200-202

HUST

#### 应该学生物信息学的N个理由



- □ 0. 计算是21世纪生物学研究的核心技能
- □ 1. 计算技能是高度可转移的
- □ 2. 计算能够帮助提高你的核心科学技能
- □ 3. 应当在博士或博后期间获得新的技能
- □ 4. 能够在生物学里建立更独特的技能
- □ 5. 可以发表更多的论文
- □ 6. 研究能有更大的灵活性
- □ 7. 工作场所不受限制
- □ 8. 计算研究的性价比高
- □ 9. 成功的科学家死在办公室
- □ 10. 能够知道为啥这个列表是从0开始



### 现代分子遗传学



- □ 分子遗传学的三大学派
  - ✿ 信息学派、结构学派和功能学派
- □ Erwin Schrödinger, 1944年, What Is Life?
  - ◆ 从信息学的角度提出遗传密码的概念,认为基因是"非周期性晶体"(大分子),生物功能的基础是分子而不是原子
  - 从量子力学的角度论证了基因的持久性和遗传模式长期稳定的可能性
  - ◆ 提出了生命"以负熵为生",从环境中抽取"序"来 维持系统的组织的概念 – 生命的热力学基础

#### WHAT IS LIFE?

The Physical Aspect of the Living Cell

ERWIN SCHRÖDINGER

### 噬菌体教堂的三主教



- Max Delbrück, Salvador Luria, Alfred
  - Hershey
    - ➡ 利用噬菌体研究基因的功能
    - ❖ 分子水平的孟德尔遗传学说
    - ◆ 1969年诺贝尔奖
- Max Delbrück (1906-1981)
  - 1930s, Niels Bohr
  - ◆ 1943, <u>fluctuation test (波动实验)</u>
  - **1944**, What is life?



#### 信息学派的发展



- □ 量子力学/量子生物学/电子生物学
  - Bernard Pullman (1919~1996)
  - **韓立International Academy of Quantum Molecular**Science
- □ 分子动力学
  - ❖ 整体牛顿力学,局部量子电动力学



# NEWS BUREAU MEMBERS HISTORY STATUTES AWARDS

CONGRESS

#### TOP NEWS

The International Academy of Quantum Molecular Sciences congratulates Martin Karplus (Univ. of Strasbourg and Harvard Univ.), Michael Levitt (Stanford Univ.), and Arieh Warshel (Univ. of Southern California) on winning the 2013 Nobel Prize in Chemistry for their work on Multiscale Models for Complex Chemical Systems.

Josef Michl President

Bioinformatics, 2025,

# The Nobel Prize in Chemistry 2013



Photo: A. Mahmoud

Martin Karplus

Prize share: 1/3



Photo: A. Mahmoud **Michael Levitt Prize share:** 1/3



Photo: A. Mahmoud

Arieh Warshel

Prize share: 1/3

### 中国早期的理论生命科学研究



- □ 用食物中的aa频率与人/小鼠的aa频率来评价食物的营养价值
- □ 输入负熵,多余的aa会被转化为其它aa
- □ 小鼠的结果未必能完全应用到人

第2卷 第1期 1962年3月 生物化学与生物物理学报

ACTA BIOCHIMICA et BIOPHYSICA SINICA

Vol. 2, No. 1 March, 1962

DIOITHOTHIAGOS, 2020, 1105/

生物体 負熵 輸入的計算 (以蛋白质营养問題为例)

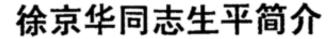
徐京華

(中国科学院生物化学研究所,上海)

表 8 各种蛋白質对青年成人信息輸入与 生物價值\*的比較

		淨 信 息	生物价值*
鸡	蛋	1.6171	96
牛	奶	1.1793	90
酵	母	0.8773	79
牛	肉	0.8114	76
明	胶	0.2131	25

<sup>\*</sup> 实驗数据取自 B. L. Oser[8]。





中国共产党党员, 1936 年参加革命, 原中国科学院上海生物化学研究所研究员、享局级离休干部徐京华同志因病长期医治无效,于 2003 年 1 月 23 日 18 时 30 分不幸逝世,享年 80 岁。

徐京华,出生于 1922 年 9 月,上海市人。1936 年 9 月参加中华民族解放 先锋队,1937 年 5 月参加中国共产党,长期从事党的地下工作。1938 年在桂林 读高中,1944 年毕业于国立西南联合大学(昆明)理学院化学系。1945 年至 1950 年担任国立中央研究院医学研究所助理员,从事神经肌肉生理研究。

解放后,徐京华同志先后担任中国科学院上海生理生化研究所助理研究员、副研究员;中国科学院上海分院生物物理研究所筹备处研究员;中国科学院上海生物化学研究所研究员、研究室主任、课题组长。先后兼任中国科学院生物物理所、理论物理所、浙江大学生物工程系教授。曾担任中国生物物理学会副理事长、上海生物物理学会理事长、上海分子科学学会理事长、上海非线性动力研究学会理事长等职。

#### DNA测序

- □ 1968年,吴瑞(Ray Wu) 发明第一个DNA测序方法
  - "引物延伸法" (Primer extension methods)
  - ★ 首次测定噬菌体λ粘端的 12个碱基



Victor Ling 林重庆



洪国藩

#### **REVIEW**

doi:10.1038/nature242

DNA sequencing at 40: past, present and future

Nations 0047 FEO 04F 0F0

Nature, 2017, 550, 345-353

#### **Technical milestones**

1953: Sequencing of insulin protein<sup>2</sup> 1965: Sequencing of alanine tRNA<sup>4</sup>

1968: Sequencing of cohesive ends of phage lambda DNA<sup>6</sup>

1977: Maxam–Gilbert sequencing<sup>9</sup>

1977: Sanger sequencing<sup>8</sup>

1981: Messing's M13 phage vector<sup>12</sup>

1986–1987: Fluorescent detection in electrophoretic sequencing<sup>14,15,17</sup>

1987: Sequenase<sup>18</sup>

1988: Early example of sequencing by stepwise dNTP incorporation 139

1990: Paired-end sequencing<sup>23</sup>

1992: Bodipy dyes<sup>140</sup>

1993: *In vitro* RNA colonies<sup>37</sup> 1996: Pyrosequencing<sup>44</sup>

1999: In vitro DNA colonies in gels<sup>38</sup>

2000: Massively parallel signature sequencing by ligation<sup>47</sup>

2003: Emulsion PCR to generate in vitro DNA colonies on beads<sup>42</sup>

2003: Single-molecule massively parallel sequencing-by-synthesis<sup>33,34</sup>

2003: Zero-mode waveguides for single-molecule analysis<sup>57</sup>

2003: Sequencing by synthesis of *in vitro* DNA colonies in gels<sup>49</sup>

2005: Four-colour reversible terminators<sup>51–53</sup>

2005: Sequencing by ligation of in vitro DNA colonies on beads<sup>41</sup>

2007: Large-scale targeted sequence capture 93-96

2010: Direct detection of DNA methylation during single-molecule sequencing<sup>65</sup>

2010: Single-base resolution electron tunnelling through a solidstate detector<sup>141</sup>

2011: Semiconductor sequencing by proton detection 142

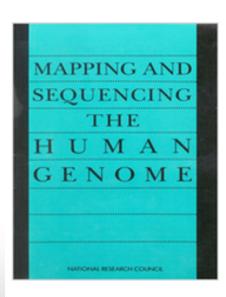
2012: Reduction to practice of nanopore sequencing 143,144

2012: Single-stranded library preparation method for ancient DNA<sup>145</sup>

### 人类基因组计划



- □ 20年, \$200,000,000/年
- □ 人类基因组测序促进人类生物学和医学的发展
- □ 测序方法的发展是应用
- □ 信息和材料的收集、分析以及分配
- □ 比较遗传学方法:解释人类基因组的信息



#### Mapping and Sequencing the Human Genome

National Research Council (US) Committee on Mapping and Sequencing the Human Genome.

Washington (DC): National Academies Press (US); 1988. ISBN-10: 0-309-03840-5

Copyright and Permissions

Hardcopy Version at National Academies Press

Search this book

### 生物信息学的发展历程



- □ 什么是生物信息学?
- ☐ Bioinformatics is a new subject of genetic data collection, analysis and dissemination to the research community (1987)
- Bioinformatics refers to database-like activities, involving persistent sets of data that are maintained in a consistent state over essentially indefinite periods of time (1994)



Hwa A. Lim

林华安

### 什么是生物信息学? (2)



□ Bioinformatics is the field of science in which biology, computer science, and information technology merge into a single discipline. The ultimate goal of the field is to enable the discovery of new biological insights as well as to create a global perspective from which unifying principles in biology can be discerned

Biology in the 21st century is being transformed from a purely lab-based science to an information science as well

from NCBI's science primer

https://www.ncbi.nlm.nih.gov/Class/MLACourse/Modules/MolBioReview/bioinformatics.html

### 广义生物信息学观点



- ☐ Biology may be viewed as the study of transmission of information: from mother cell to daughter cell, from one cell or tissue type to another, from one generation to the next, and from one species to another. This informational viewpoint is termed bioinformatics
- □ 生物学研究可以被看成是研究信息的传递: 从 DNA经转录翻译到蛋白质,从细胞质中到细胞核内,从母细胞到子细胞,从一个细胞或一个组织到另一个细胞或另一个组织,从一代到下一代,从一个物种到另一个物种的进化演变。这种信息论的观点可称为生物信息学 (Eisenberg et al., 2006)

### 生物信息学发展过程中的里程碑性事件



时间	重大事件
1962	鲍林的分子进化理论
1965	Margaret Dayhoff 构建蛋白质序列图谱
1970	Needleman-Wunsch 算法开发
1977	DNA 测序以及使用利用 Staden 软件进行分析
1981	Smith-Waterman 算法开发
1981	序列模体的概念(motif)
1983	序列数据库搜索算法的开发(Wilbur-Lipman)
1985	快速序列相似性搜索工具 FASTP/FASTN 的开发
1990	快速序列相似性搜索工具 BLAST 的开发
1996	酵母基因组的完全测序
1997	PSI-BLAST 工具的开发
1998	秀丽线虫基因组的完全测序
1999	果蝇基因组的完全测序

# 基因组、计算和应用的里程碑性事件



•					
Ger	10me	mı	lest	on	es

1977: Bacteriophage ΦΧ174 (ref. 72) 1981: Smith–Waterman<sup>156</sup>

1982: Bacteriophage lambda<sup>13</sup>

1995: Haemophilus influenzae<sup>26</sup>

1996: Saccharomyces cerevisiae<sup>27</sup>

1998: Caenorhabditis elegans<sup>28</sup>

2000: Drosophila melanogaster<sup>32</sup>

2000: Arabidopsis thaliana 146

2001: Homo sapiens<sup>29–31</sup>

2002: Mus musculus<sup>147</sup>

2004: Rattus norvegicus<sup>148</sup>

2005: Pan troglodytes<sup>149</sup>

2005: Oryza sativa<sup>150</sup>

2007: Cyanidioschyzon merolae<sup>126</sup>

2009: Zea mays<sup>151</sup>

2010: Neanderthal<sup>88</sup>

2012: Denisovan<sup>145</sup>

2013: The HeLa cell line 152,153

2013: Danio rerio<sup>154</sup>

2017: Xenopus laevis<sup>155</sup>

#### Computational milestones

1982: GenBank (https://www.nct

1990: BLAST16

1995: TIGR assembler<sup>24</sup>

1996: RepeatMasker

1997: GENSCAN157

1998: phred, phrap, consed<sup>22</sup>

2000: Celera assembler<sup>25</sup>

2001: Bioconductor

2001: EULER<sup>74</sup>

2002: BLAT<sup>158</sup>

2002: UCSC Genome Browser<sup>159</sup>

2002: Ensembl<sup>160</sup> 2005: Galaxy<sup>161</sup>

2007: NCBI Short Read Archive

2008: ALLPATHS<sup>162</sup>

2008: Velvet<sup>75</sup>

2009: Bowtie<sup>83</sup>

2009: BWA82

2009: SAMtools84

2009: BreakDancer<sup>163</sup>

2009: Pindel<sup>164</sup>

2009: TopHat<sup>115</sup>

2010: SOAPdenovo165

2010: GATK<sup>85</sup>

2010: Cufflinks<sup>116</sup>

2011: Integrated Genomics Viewer<sup>166</sup>

2013: HGAP/Quiver<sup>167</sup>

2017: Canu<sup>81</sup>

#### **Application milestones**

1977: Genome sequencing<sup>72</sup>

1982: Shotgun sequencing<sup>13</sup>

1983, 1991: Expressed sequence tags 107,108

1995: Serial analysis of gene expression 109

1998: Large-scale human SNP discovery<sup>168</sup>

2004: Metagenome assembly 122

2005: Bacterial genome resequencing with NGS<sup>40,41</sup>

2007: Chromatin immunoprecipitation followed by sequencing

(ChIP-seq) using NGS<sup>117</sup>

2007–2008: Human genome and cancer genome resequencing using

NGS<sup>55,90-92</sup>

2008: RNA-seq using NGS<sup>110–114</sup>

2008: Chromatin accessibility using NGS<sup>118</sup>

2009: Exome resequencing using NGS<sup>97</sup>

2009: Ribosome profiling using NGS<sup>119</sup>

2010: Completion of Phase I of the 1000 Genomes Project<sup>98</sup>

2010: De novo assembly of a large genome from short reads 169

2011: Haplotype-resolved human genome resequencing using NGS<sup>170,171</sup>

2016: Human genome *de novo* assembly with PacBio<sup>172</sup>

2017: Human genome *de novo* assembly with nanopore<sup>64</sup>

025, HUST

### 分子进化理论: 鲍林



#### Molecular Disease and Evolution

by Linus Pauling

(Rudolf Virchow Lecture, 5 November 1962)



1901 ~ 1994

I believe that it is likely that a human being manufactures 50,000 or 100,000 different kinds of protein molecules. A representative protein molecule, such as hemoglobin, is built of about 10,000 atoms. It has a well defined structure; for most of the protein molecules not a single atom is out of place.

### Dayhoff 打分矩阵



- □ 1978年, 34个蛋白质超家族
- □ 序列相似性≥85%
- □ 1572个点突变
- □ PAM: Accepted Point Mutation,可 接受的点突变
- □ 马尔可夫模型 (Markov Model): 位点 突变速率独立均等
- □ PAM1: 序列分歧~1%时的氨基酸替代 打分矩阵



Margaret Dayhoff, 1925 ~ 1983

### 生物信息学: 爸爸/妈妈是哪位?





弗雷德里克·桑格 是我?



莱纳斯·卡尔·鲍林 不是我?



坦布尔·斯密 右边说是我



乔治·贝尔 左边这位!



鲍琳•霍奇维格 我都说是我了



**薛定谔** 关我啥事啊?



林华安 其实应该是我



玛格丽特·奥克雷·戴霍夫 那我呢?

Bioinformatics, 2025, HUST

### 国际四大核酸序列数据库



- □ 1974年,George I. Bell等人收集DNA序列,构建 GenBank数据库,1982开发第一个版本
- □ 1980年, 欧洲EBI-ENA数据库建立
- □ 1984年,日本DDBJ数据库建立成立
- □ 2016年,中科院北京基因组所BIG Data Center at Beijing Institute of Genomics(BIGD)数据库建立;中科院上海营养与健康研究所Bio-Med Big Data Center (BMDC)数据库建立
- □ 2019年, BIGD改为NGDC (National Genomics Data Center); BMDC改成NGDC/NODE (National Omics Data Encyclopedia)

#### **NGDC**



□ <a href="https://ngdc.cncb.ac.cn/">https://ngdc.cncb.ac.cn/</a>



### 2019新型冠状病毒信息库



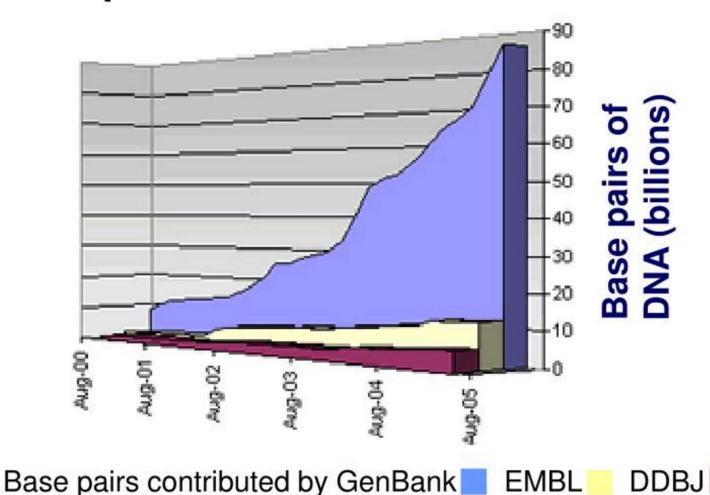
□ <a href="https://ngdc.cncb.ac.cn/ncov/">https://ngdc.cncb.ac.cn/ncov/</a>



### 核酸数据库的数据增长



#### Growth of the International Nucleotide Sequence Database Collaboration



### 获取序列及检索公共数据库



- □ Entrez的开发, David J. Lipman等人
- □ 提供关键字的搜索的方法
- □ "硬搜索":包含关键字的,完全匹配的结果
- □ "软搜索":与查询内容相关的信息
- □ 查询内容: 基因/蛋白质的名称、标识符, 文献、

蛋白质结构等

#### Search NCBI



#### https://www.ncbi.nlm.nih.gov/search/



#### **COVID-19 resources**



#### **COVID-19 treatment**

Coronavirus Disease 2019 (COVID-19) Treatment Guidelines (National Institutes of Health)



#### SARS-CoV-2 genomes

Download viral genome and protein sequences, annotation and a data report



#### **NCBI Virus**

The most up-to-date set of SARS-CoV-2 nucleotide and protein sequences

COVID-19 news >

#### Literature

#### PubMed

PubMed® comprises more than 35 million citations for biomedical literature from MEDLINE, life science journals, and online books. Citations may include links to full text content from PubMed Central and publisher web sites.



**Example searches** Search for titles, citations, identifiers and more

Cancer Statistics 2021

Tunyasuvunakool, Nature 2021

Revealing protein-protein interactions by transcriptome sequencing

GeneReviews

#### Literature databases

#### Bookshelf

Books and reports

#### MeSH

Ontology used for PubMed indexing

#### **NLM Catalog**

Books, journals and more in the NLM Collections

#### PubMed

Scientific and medical abstracts/citations

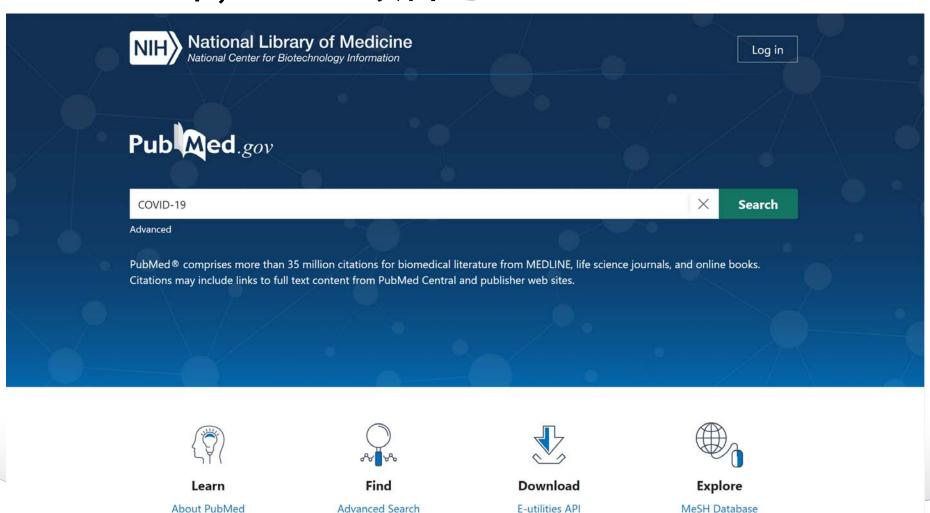
#### **PubMed Central**

### PubMed: 文献检索



#### □ 1992年, Entrez项目之一

FAQs & User Guide



FTP

Journals

Clinical Queries

# PubMed: 文献检索



National Library of Medicine  National Center for Biotechnology Information  Log in				
Pub Med.gov	COVID-19  Advanced Create alert Create RSS	X Search User Guide		
	Save Email Send to	Sorted by: Best match Display options 🗱		
my ncbi filters 🖪	332,931 results	<pre></pre>		
RESULTS BY YEAR	Treatment Mechanism Transmission	ed Clinical Queries to refine your search  More filters  sequence, and clinical content from NCBI		
1981 2023	COVID-19 diagnosis -A review o  Yüce M, Filiztekin E, Özkaya KG.	f current methods.		
Abstract  Free full text  Full text	PMID: 33126180 Free PMC article.  Share A fast and accurate self-testing tool for Co	52. doi: 10.1016/j.bios.2020.112752. Epub 2020 Oct 24. Review.  OVID-19 diagnosis has become a prerequisite to comprehend to take medical and governmental actions accordingly. SARS-		

### 精准医学



Next >

# □ 2011年,新的疾病分类方法 (New Taxonomy)



#### **Toward Precision Medicine**

Building a Knowledge Network for Biomedical Research and a New Taxonomy of Disease

National Research Council (US) Committee on A Framework for Developing a New Taxonomy of Disease.

Washington (DC): National Academies Press (US); 2011. ISBN-13: 978-0-309-22222-8 ISBN-10: 0-309-22222-2

Copyright and Permissions

Hardcopy Version at National Academies Press

#### Excerpt

A Knowledge Network of Disease could embrace and inform rapidly expanding efforts by the biomedical research community to define at the molecular level the disease predispositions and pathogenic processes occurring in individuals. This network has the potential to play a critical role across the globe for the public-health and health-care-delivery communities by enabling development of a more accurate, molecularly-informed taxonomy of disease.

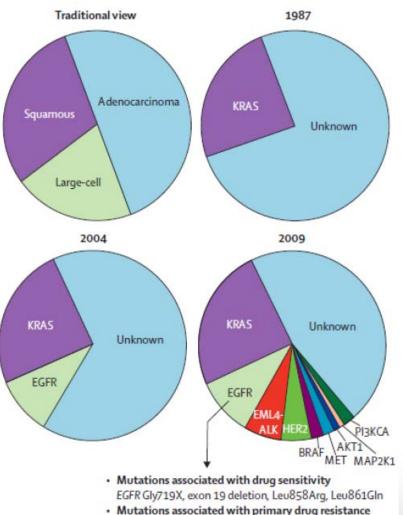
http://www.ncbi.nlm.nih.gov/books/NBK91503/

### 新的疾病分类方法



#### □分子组学数据

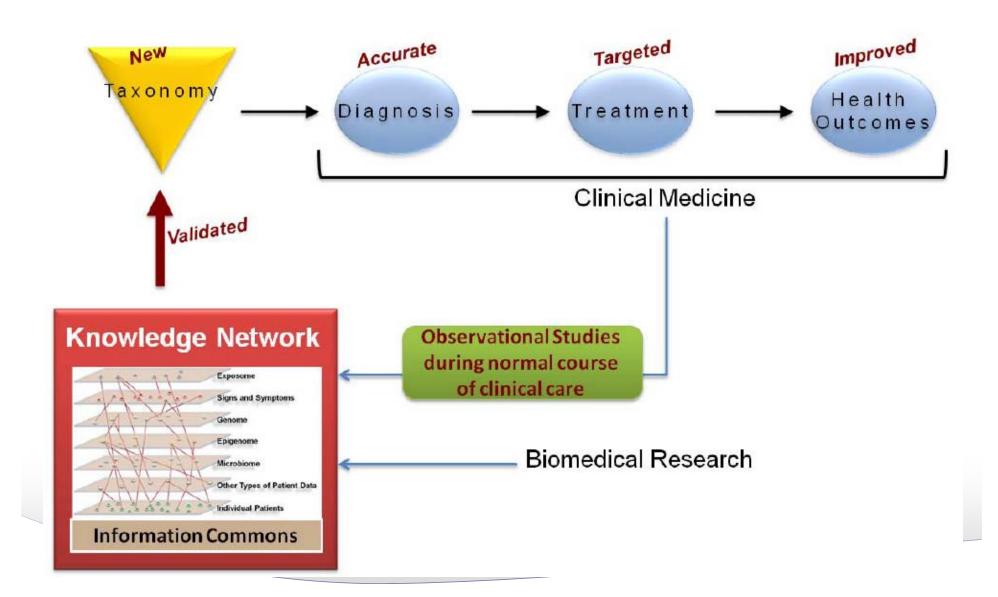
- ✿ Genome (基因组)
- ♣ Transcriptome (转录组)
- ♣ Proteome (蛋白质组)
- ◆ Metabolome (代谢组)
- ♣ Lipidome (脂质组)
- ♣ Epigenome (表观组)
- ✿ Meta-genome(宏基因组)



- Mutations associated with primary drug resistance EGFR exon 20 insertions
- Mutations associated with acquired drug resistance EGFR Thr790Met, Asp761Tyr, Leu747Ser, Thr854Ala

## 信息共享 (Information Commons)



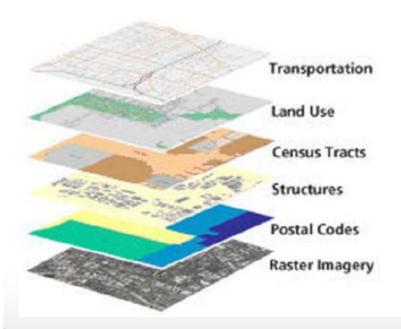


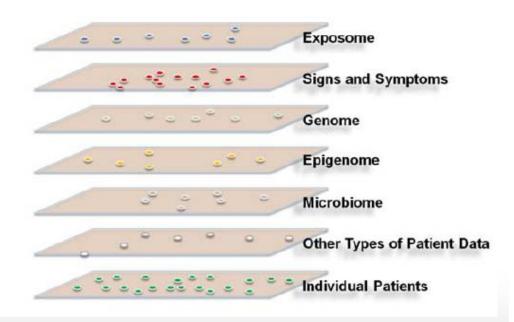
# 信息共享 (Information Commons)



Google Maps: GIS layers
Organized by Geographical Positioning

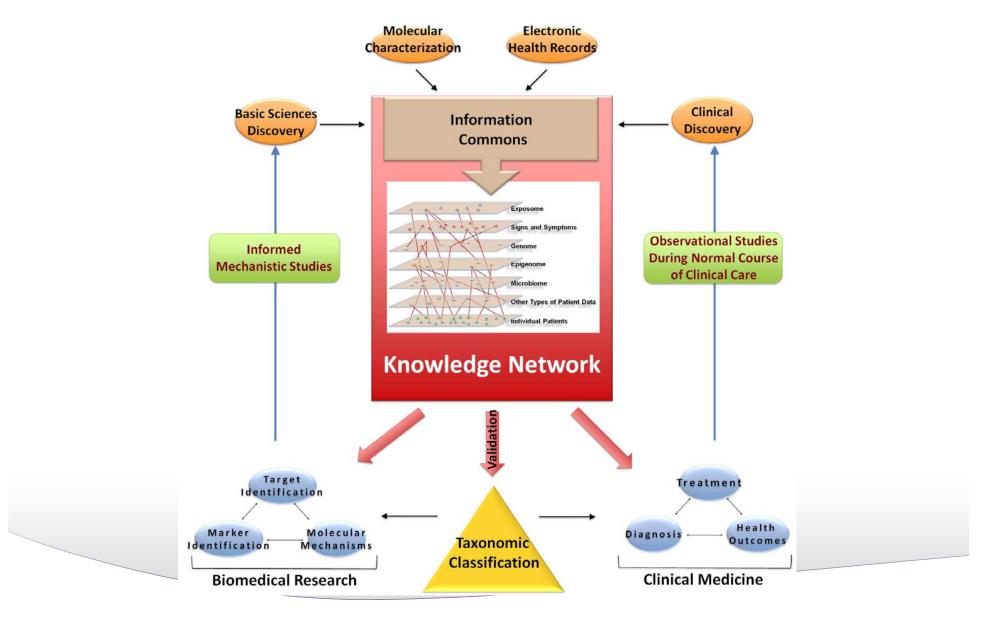
### Information Commons Organized Around Individual Patients





### **Biomedical Knowledge Network**





### 生物信息学主要研究内容与方向



- □ 开发新的算法及统计学方法来揭示<u>大数据</u> 之间的联系
- □ 开发、设计相关的数据库和工具,能够方便有效的获取、管理以及使用各种类型的数据和信息
- □ 分析和解释各种类型的生物学数据,包括 核酸、氨基酸序列、蛋白质功能结构域以 及蛋白质三级结构等

http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/milestones.html

# 中国生物信息学的研究方向



1. 基因组信息学	12. 药物信息学
2. 转录组信息学	13. 人工生物系统的设计与控制
3. 表观组信息学	14. 生物信息数据挖掘方法
4. 蛋白质组信息学	15. 生命与健康大数据科学
5. 修饰组信息学	16. 精准健康信息学
6. 结构信息学	17. 群体遗传与计算演化基因组学
7. 单细胞生物信息学	18. 生物影像信息学
8. 三维基因组信息学	19. 植物信息学
9. 微生物组信息学	20. 生物信息学研究新方向、新技
10. 计算癌症基因组学	术与新方法
11. 计算系统生物学	21. 人工智能与生物学

### 中国生信相关的学会、协会(二级)



中国生物信息学学会(筹) 中国细胞生物学学会功能基因组信息学与系统生物学分会 中国生物工程学会计算生物学与生物信息学专业委员会 中国计算机学会生物信息学专委会 中国遗传学会基因组学专业委员会 中国生物物理学会生物信息学与理论生物物理专业委员会 中国生物化学与分子生物学会分子系统生物学专业委员会 中国运筹学会计算系统生物学分会 中国人工智能学会生物信息学与人工生命专业委员会 中国电子学会生物计算与生物信息处理专业委员会 中国医药生物技术协会生物医学信息技术分会 中国交叉科学学会生物信息学专业委员会 中国系统仿真学会生命系统建模仿真专业委员会 中国生物物理学会人工智能生物学分会

### 湖北省生物信息学会



- □ 2018年10月16日,湖北省科学技术协会批准
- □ 2018年12月8日,湖北省生物信息学会成立大会、第四届湖北省生物信息学青年学者论坛暨在汉四校两所联合生物信息研究生年会
- □ 2019年3月29-31日,第六届华中-华东地区生物信息学研讨会,宜昌,三峡大学
  - ♥ 终身成就奖 & 杰出贡献奖
- □ 征集学会的logo
- □ 生物信息分析认证体系: 等级考试



## 国内相关专业介绍



- □ 华中科技大学
  - 生命学院: 宁康、郭安源、薛宇、陈卫华
  - 物理学院: 肖奕、黄胜友、刘士勇
  - ➡ 其他学院:潘林强、何西淼、王超龙、夏天
- □ 武汉大学:刘娟、赵华斌、周宇、邹秀芬
- □ 华中农业大学: 张红雨、陈玲玲、李国亮、谢为博
- □ 北京大学:魏丽萍、高歌、李川昀、来鲁华、欧阳颀、邓明华、席瑞斌、崔庆华、陆剑、李程
- □ 清华大学: 张学工、汪小我、李梢、鲁志、杨雪瑞、张强锋、 曾坚阳
- □ 上海同济大学:刘小乐、江赐忠、张勇、曹志伟、刘琦

### 985、211高校



- □ 厦门大学:纪志梁
- □ 中国科学技术大学: 刘海燕、李骜、瞿昆
- □ 复旦大学:赵兴明、倪挺、田卫东
- □ 天津大学:高峰
- □ 浙江大学:陈铭、樊龙江
- □ 上海交通大学:沈红斌、王侃侃
- □ 中山大学: 贺雄雷、任间、谢志、骆观正、王金凯
- □ 南方医科大学:张镇海、王栋
- □ 四川大学:沈百荣、谢丹
- □ 中国农业大学: 王向峰

## 中科院系统 & 研究所



□ 北京中科院生物物理所: 陈润生 □ 北京基因组所:章张、王前飞、张治华、方向东 北京中科院遗传所:王秀杰、钱文峰 □ 北京中科院计算所:赵屹 □ 北京中科院数学与系统科学研究院: 张世华、王勇 □ 北京生命科学研究院:赵方庆、吴金雨 □ 北京动物所:张勇 □ 上海生物信息中心: 李亦学、谢鹭 □ 上海系统生物学实验室: 陈洛南 □ 上海马普所:王泽峰、韩敬东、徐书华、杨力、李海鹏

□ 中国医学科学院基础医学研究所: 蒋太交

# 知名企业: 华大基因







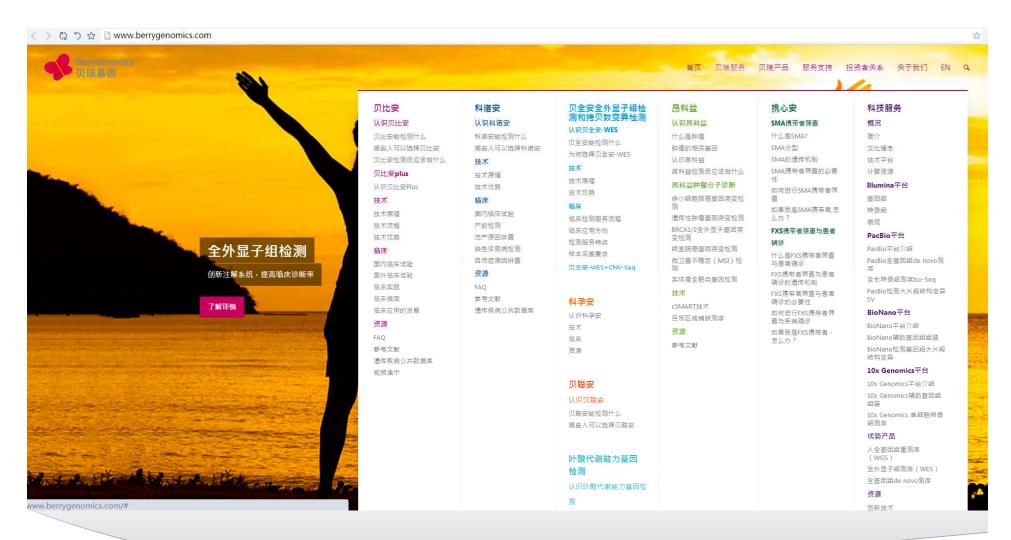
# 知名企业: 诺禾致源





# 知名企业: 贝瑞基因





# 知名企业: 武汉未来组





#### F/J == / /



#### 动植物基因组 de novo 3.0

重新定义大基因组组装指标 Contig N50>1M Scaffold N50>5M



#### 全长转录组

无需拼接的全长转录本 精确鉴定可变剪切和融合基因信息



#### 细菌基因组 完成图3.0

一个重叠群一染色体 无Gap,无N碱基错误

### **Bioinformatics in China**



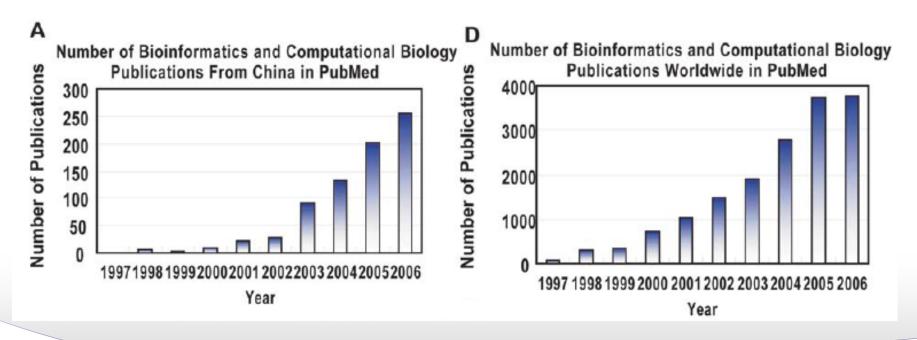
OPEN & ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

#### Perspective

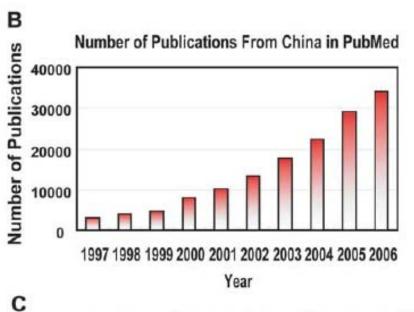
### Bioinformatics in China: A Personal Perspective

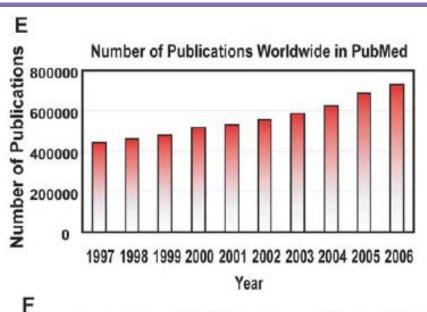
Liping Wei<sup>1</sup>\*, Jun Yu<sup>2</sup>\*

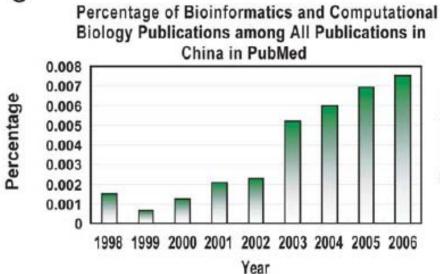


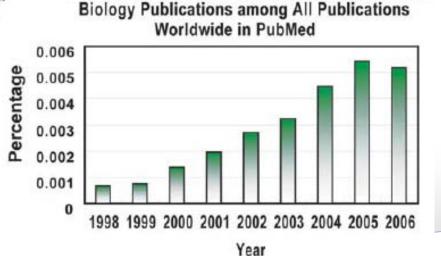
### **Bioinformatics in China**











Percentage of Bioinformatics and Computational

## 相关知识储备



- □ 生物学背景: e.g., 细胞生物学、分子生物学、发育生物学、分子神经生物学、生物化学, ...
- □ 机器学习/深度学习/生成式AI
- □ 计算能力/编程能力: Perl/Python, R, PHP+MySQL, JAVA...
- □ 分子进化理论: MP, NJ, ML...
- □ 生物统计学
- □ 结构生物学

## 生物信息学的相关杂志



#### 生物信息学相关期刊名称

#### 网址

Bioinformatics http://bioinformatics.oxfordjournals.org/

BMC Bioinformatics http://www.biomedcentral.com/bmcbioinformatics/

Genome Biology http://genomebiology.com/

Genome Research http://www.genome.org/

Nucleic Acids Research http://nar.oxfordjournals.org/

Briefings in Bioinformatics http://www.henrystewart.com/briefings\_in\_bioinformatics/

FEBs letters http://www.febsletters.org/

Biochemical and Biophysical Research http://www.sciencedirect.com/science/journal/0006291X

Communications

Molecular Systems Biology http://www.nature.com/msb/index.html

Molecular Biology and Evolution http://mbe.oxfordjournals.org/

PLoS Computational Biology http://www.ploscompbiol.org/

PLoS ONE http://www.plosone.org/

Protein Science http://www.proteinscience.org/

Proteins http://www3.interscience.wiley.com/cgi-bin/jhome/36176

Protein Engineering Design and Selection http://peds.oxfordjournals.org/

### **GPB: Bioinformatics Commons**





Volume: 16. Issue: 4



Volume: 16, Issue: 5

#### Preface

Bioinformatics Commons: The Cornerstone of Life and Health Sciences

Zhang Zhang, Yu Xue, Fangging Zhao

View abstract

Page 223-225

Download **⊕** 196

#### Database

CIRCpedia v2: An Updated Database for Comprehensive Circular RNA Annotation and Expression Comparison

Rui Dong, Xu-Kai Ma, Guo-Wei Li, Li Yang

View abstract

Page 226-233

#### Editorial

A Scientist Guerilla Fighter in the Frontiers of Bioinformatics—In Memory of Bailin Hao

Jun Yu

View abstract

Page 307-309

Download 39

#### **Original Research**

Polyphyly in 16S rRNA-based LVTree Versus Monophyly in Whole-genome-based CVTree

Guanghong Zuo, Ji Qi, Bailin Hao

View abstract

Page 310-319

## 如何评价生物信息学研究的水平?

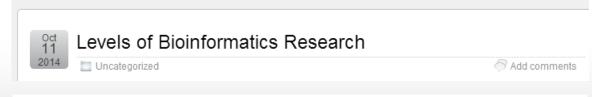


- □ 0级 (Level 0): 为建模、而建模(modeling for modeling's sake)
- □ 1级(Level 1, 菜鸟级): 给数据、能分析
- □ 2级(Level 2, 肉鸟级): 想新招、玩数据
- □ 3级(Level 3, 顶级): 玩数据、作发现
- □ X级(Level X, 神级): 玩科学、讲政治



刘小乐教授 哈佛大学

#### 合适的研究体系、好的工具、百折不挠的毅力!



http://www.longwoodgenomics.org/2014/10/11/levels-of-bioinformatics-research/

## Homolog.us的评价体系



- Layer 1 Using web to analyze biological data
- □ Layer 2 Ability to install and run new programs
- Layer 3 Writing own scripts for analysis in PERL, python or R
- Layer 4 High level coding in C/C++/Java for implementing existing algorithms or modifying existing codes for new functionality
- □ Layer 5 Thinking mathematically, developing own algorithms and implementing in C/C++/Java

Layer 1-4 = Level 1; Layer 5 = Level 2

A beginner's guide to bioinformatics - part I A beginner's guide to bioinformatics - part II

http://www.homolog.us/blogs/blog/2011/07/2 2/a-beginners-guide-to-bioinformatics-part-i/

http://www.homolog.us/blogs/blog/2011/07/2 2/a-beginners-guide-to-bioinformatics-part-ii/

# 生物信息学要不要做实验?



- □ 生物信息学家应该做实验
- □ 最新的技术, Computational biologists should stay at the cutting edge of technologies
- □ 高通量
- □ 关注新技术产生的数据的质量
- □ 要回答具体的生物学问题
- □ 实验不能太难,也不能太贵