



面向世界科技前沿，面向国家重大需求，面向国民经济主战场，率先实现科学技术跨越发展，率先建成国家创新人才高地，率先建成国家高水平科技智库，率先建设国际一流科研机构。

——中国科学院办院方针



## 北京基因组所等开发完成基于K-mer的基因组组分析数据库

文章来源：北京基因组研究所 发布时间：2015-11-02 【字号：小 中 大】

我要分享

在过去的几十年中，人们往往使用高度保守的基因家族进行系统进化分析，采用全基因组序列进行系统进化分析并不普遍。目前，基于是否进行序列的比对，分子系统发生树的构建分为两类。其中，不需要进行序列比对的方法是依据K-mer向量计算的距离矩阵进行系统进化分析，大量的研究证实该算法是行之有效的，尤其是对基因组中诸如蛋白编码序列等的特定区域。不仅如此，K-mer算法还在组学的其他方面，包括基因组组装、motif预测、重复序列的识别以及基因组的复杂性评估等都受到了广泛的关注。基于K-mer算法在组学中的重要表现，在这个大规模基因组数据快速积累的时代，构建一个基于K-mer算法易于存储并且将大量基因组数据可视化处理的数据库十分迫切。

为此，中国科学院北京基因组研究所基因组科学与信息重点实验室于军组和英国伦敦大学学院（UCL）肿瘤研究所王大鹏合作开发了一套基于K-mer算法的基因组组分析数据库KGCAK。此项研究于近期发表在Biology Direct 杂志。

在这个数据库中，研究人员搜集了Ensembl、Phytozome和NCBI等几大主流基因组数据库中包括高等动植物、原生物、真菌、细菌、病毒等在内的8000多个核基因组或者细胞器基因组，同时包括基因组不同维度的序列，主要有DNA、cDNA、CDS、氨基酸和ncRNA序列。并且还分别计算和存储了核酸序列（K从2变化到10）和氨基酸序列（K从1变化到5）的K-mer向量，以方便进行不同维度数据跨物种的系统发生树构建。此外，该数据库提供了评估不同物种基因组复杂度的交互工具，主要包括基因组基本特征参数、K-mer向量的数学参数统计、频率分布、唯一性比率，以及二维和三维空间可视化分析基因组参数和K-mer参数的交互关系等。

总的来说，该数据库通过捕获基因组序列特征并把基因组转化成更易于理解和可视化的数字K-mer向量，以期通过K-mer算法用可视化的图形和定量的数据构建一个比较基因组学的平台，将为系统发生树构建和通过基因组数据研究物种关系提供良好的参照和指引。

文章链接

### 热点新闻

#### 中科院与广东省签署合作协议 ...

发展中国家科学院中国院士和学者代表座...  
白春礼在第十三届健康与发展中山论坛上...  
中科院江西产业技术创新与育成中心揭牌  
中科院西安科学园暨西安科学城开工建设  
中科院与香港特区政府签署备忘录

### 视频推荐

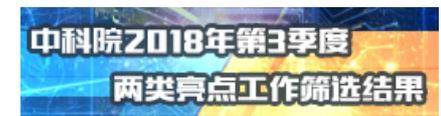


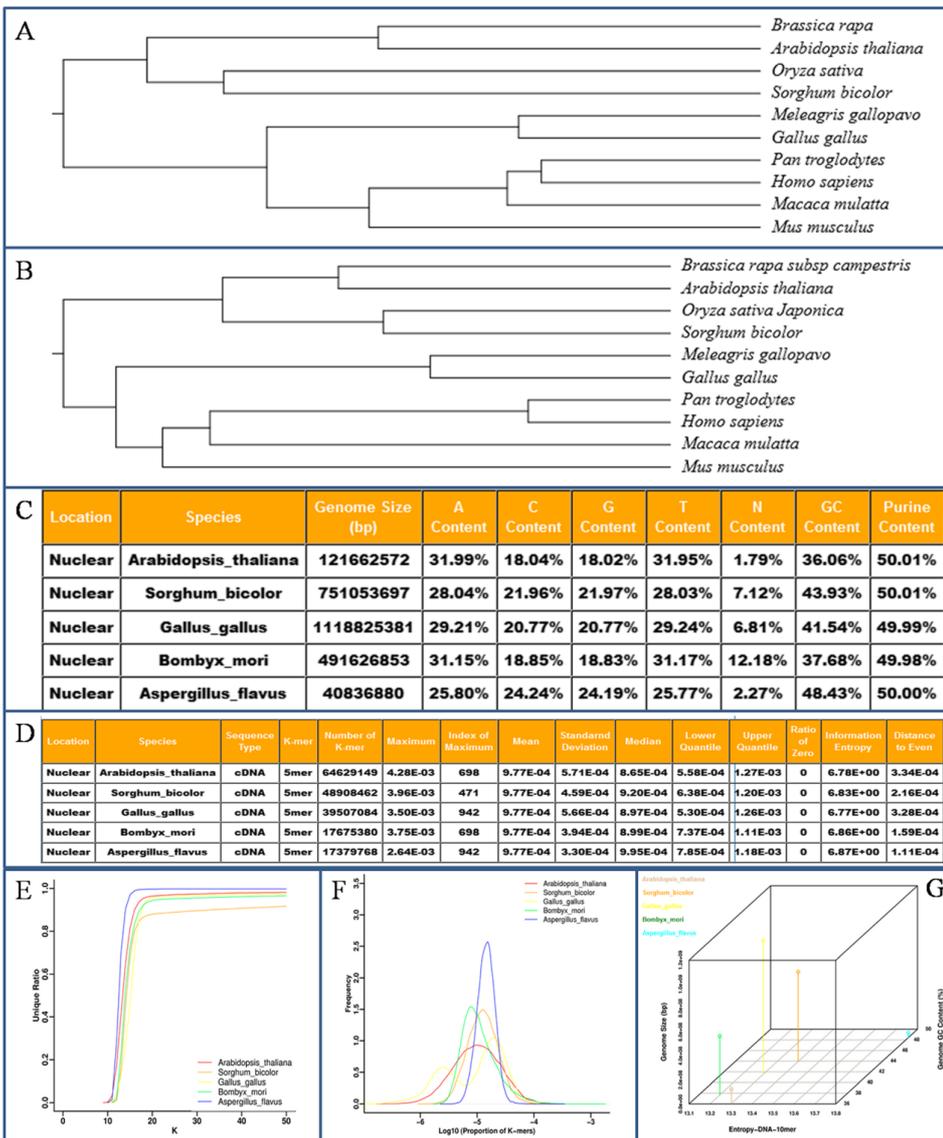
【新闻联播】“率先行动”计划 领跑科技体制改革



【新闻直播间】我科学家造血干细胞研究获新突破

### 专题推荐





KGCAK数据库中基本功能模块举例

(责任编辑：叶瑞优)



© 1996 - 2018 中国科学院 版权所有 京ICP备05002857号 京公网安备110402500047号 联系我们

地址：北京市三里河路52号 邮编：100864