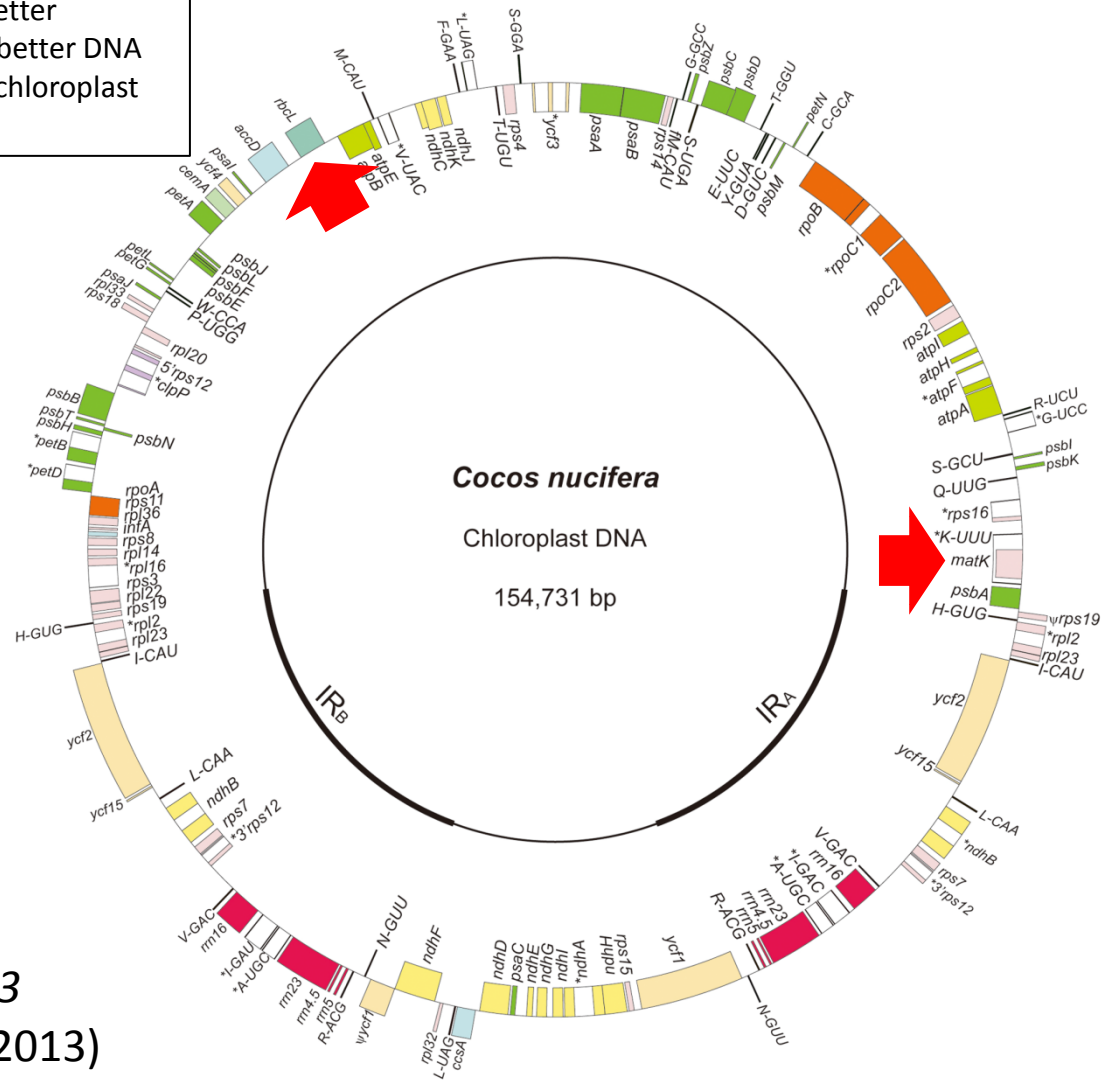2015-11-26

# 3. DNA Barcode

- **Plant barcodes(matK, rbcL)**
- **BOLD/GBIF DB**
- **GBIF entries of 'Palm Trees'**

# Plant DNA barcodes

In 2009, a collaboration of a large group of plant DNA barcode researchers proposed two chloroplast genes**, rbcL and matK**, taken together, as a barcode for plants.[6] Adding the nuclear internal transcribed spacer ITS2 region was proposed to provide better resolution between species.[21] As of 2015, the search for better DNA barcodes for plants continues, with the proposal that the chloroplast region *ycf1* may be suitable.

https://en.wikipedia.org/wiki/DNA_barcoding



*Cocos nucifera*

Chloroplast DNA

154,731 bp

PMID:*24023703*
(Huang et al., 2013)

# GBIF (Global Biodiversity Information Facility)

http://www.gbif.jp/bol/

JBIF 地球規模生物多様性情報機構日本ノード
Japan Node of Global Biodiversity Information Facility

バーコードオブライフデータを用いた生物種同定システム　　ホーム > バーコードオブライフデータを用いた生物種同定システム

- システムの概要
  本システムはCOIやITSなどの塩基配列を バーコードオブライフデータベース(BOLD)および公共 DNAデータベース(DDBJ)から抽出してデータベースを構築しています。 任意のDNA配列の生物種 を同定することが可能なシステムです。
- 参照するデータベース
  以下のデータベースのいずれかを選択してください。

■代表配列のデータベース
BOLDおよびDDBJには同じ生物種名の塩基配列が含まれています。 生物種を同定するためにはその冗長性が 扱いづらいため、1生物種につき 1件の代表塩基配列としたデータベースを構築しています。

・BOLD由来のデータベース
- COI-5P (97,965件) データ更新日:2015年02月06日　　FASTA (15MB) リストファイル (1.8MB)
- COI-3P (4,693件) データ更新日:2015年02月06日　　FASTA (0.7MB) リストファイル (84KB)
- rbcL (34,121件) データ更新日:2015年02月06日　　FASTA (7.1MB) リストファイル (0.6MB)
- ITS (22,378件) データ更新日:2015年02月06日　　FASTA (4.4MB) リストファイル(0.4MB)
- matK (34,656件) データ更新日:2015年02月06日　　FASTA (8.8MB) リストファイル (0.6MB)

Unique sequences

・DDBJ由来のデータベース
- 16S rRNA 細菌(233,506件) データ更新日:2015年02月06日 FASTA (45MB) リストファイル (3.4MB)

■全件のデータベース
データベースの塩基配列全件に対して比較することが可能です。

・BOLD由来のデータベース
- COI-5P (2,807,009件) データ更新日:2015年02月06日　　FASTA (380MB) リストファイル
- COI-3P (20,898件) データ更新日:2015年02月06日　　FASTA (3MB) リストファイル (0
- rbcL (77,415件) データ更新日:2015年02月06日　　FASTA (13MB) リストファイル

右側メニュー:
- GBIF／JBIFとは
- GBIFデータの利用
- GBIFへのデータ登録
- 各種ドキュメント
- 関連の活動
- Barcode同定
- リンク

GBIFニュースレター(日本語版)
GBits

GBIF　Japan Node
↓
JBIF

Dr.Yamazaki(NBRP) supports JBIF database.

BOLD and DDBJ sources

## BOLD(Barcode of Life Data)

| Sequence statistics | | Species coverage (formally described) | |
|---|---|---|---|
| Barcode clusters for animals (BINs) | 382,631 | Animals | 154,900 |
| All Sequences | 4,321,441 | Plants | 58,701 |
| Barcode Sequences | 3,761,354 | Fungi & Other Life | 16,760 |

2015/2/27

# Download files: GBIF matK sequences

(1)matK_rpsv.list    (TSV format: Accession ID, Species name, barcode name, GenBank ID)

```
GBVH547-11^ Guatteria olivacea^ matK^     AY740940
GBVA1687-11^Biarum carduchorum^ matK^     EU886521
POWNA1560-12^   Kickxia spuria^ matK^     JN894552
GBVE3433-11^Raphanus sativus var. raphanistroides^  matK^    AB354261
GBVJ1159-11^Ceanothus foliosus var. vineatus^   matK^    AF049803
GBVR3836-13^Cylindropuntia cholla^  matK^    FN997446
GBVD1799-11^Carex vexans^    matK^    GU173775
GBVS4700-13^Opuntia pumila^ matK^    JF786826
```

(2)matK_rpsv.fasta   (fasta format: >Accession ID, Sequence)

```
>GBVH547-11
TACCTCACCCCGCCCATCTGGAAATCTTGGTTCAAATATTTCGCTCTTGGATACAAGATGCCCCCTCTTTGCATTTATTGCGATCCTTTC
>GBVA1687-11
TTTGCTGTCATTATGGAAATTCCTTTCTCATTGCGACTAGTATACTCCCTCGAAGAAAAAAAGAAATACCAAAATCTCAGAATTTACGA
>POWNA1560-12
TCACATTTAAATTTTGTGTTAGATATACTAATACCCTACCCTGTCCATGTGGAAATCTTGGTTCAAACTCTTCGCTATTGGGTAAAAGAT
>GBVE3433-11
ATGTGTCATTTCAGAACTCAAGAAAATAAAGACTTTACTTTTAGTTCAAATCGAATTTCAATCCAAATGGAGAAATTTCAAGGATATTTA
>GBVJ1159-11
ATGGAAGAGTTTCAAGGATATTTCGAACTAAATAGATCTCGGCAACACGATCTCCTATACCCACTTATCTTTCGGGAGTATATTTATGCA
>GBVR3836-13
```

# Reference : Palm trees and matK entries



PMID:*24023703*
(YY Huang  et al., 2013)

| Genus | Ex. Species name | matK entries |
| --- | --- | --- |
| Cocos | *Cocos nucifera*(Coconut) | 1 |
| Phoenix | *Phoenix dactylifera(Date Palm)* | 6 |
| Bismarckia | *Bismarckia nobilis* | 1 |
| Pseudophoenix | *Pseudophoenix lediniana* | 6 |
| Chamaedorea | *Chamaedorea elegans* | 31 |
| Elaeis | *Elaeis guineensis(*Oil Palm*)* | 1 |
| Calamus | Calamus sp.(Rattan) | 42 |
| Areca | *Areca catechu(*Betel nuts *)* | 2 |
| Metroxylon | *Metroxylon salomonense* | 1 |

## 2015-11-26

## 4. NGS Read Alignment

- **DDBJ Pipeline**
- **SAM/BAM format**
- **Visualization(Samtools tview)**
  - **SRA100551(query)**
  - **GU811709 (ref.)**

# Date palm: datasets of chloroplast genome

Do not download data: the next page tool imports automatically

■ Phoenix dactylifera (date palm) : taxid:42345

■ **GU811709 (reference sequence)**

http://www.ncbi.nlm.nih.gov/genome/organelles/2664?

**Phoenix dactylifera**

Items 1 - 2 of 2    << First    < Prev    Page 1    of 1    Next >

| Organism | Name | RefSeq | INSDC | Size (Kb) | GC(%) | Protein | rRNA | tRNA | Other RNA | Gene | Pseudogene |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Phoenix dactylifera | Pltd | NC_013991.2 | GU811709.2 | 158.46 | 37.2 | 95 | 8 | 44 | - | 149 | 2 |
| Phoenix dactylifera | MT | NC_016740.1 | JN375330.1 | 715 | 45.1 | 43 | 3 | 18 | - | 44 | 1 |

chloroplast genome

■ **SRA100551 (query sequences)**

9 cultivars

```
−<SAMPLE_SET>
  +<SAMPLE center_name="The University of Texas at Austin" alias="AJW" accession="SRS478070"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="PER" accession="SRS478072"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="SUK-A" accession="SRS478078"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="DEK" accession="SRS478079"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="SUK-Q" accession="SRS478080"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="RAB" accession="SRS478081"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="SHA" accession="SRS478082"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="MOS-A" accession="SRS478083"></SAMPLE>
  +<SAMPLE center_name="The University of Texas at Austin" alias="MOS-H" accession="SRS478084"></SAMPLE>
</SAMPLE_SET>
```

```
<LIBRARY_STRATEGY>WGS</LIBRARY_STRATEGY>
<LIBRARY_SOURCE>GENOMIC</LIBRARY_SOURCE>
```

whole genome sequencing

OPEN ACCESS Freely available online    PLOS ONE

Whole Mitochondrial and Plastid Genome SNP Analysis of Nine Date Palm Cultivars Reveals Plastid Heteroplasmy and Close Phylogenetic Relationships among Cultivars

Jamal S. M. Sabir[1], Dhivya Arasappan[2], Ahmed Bahieldin[1,3], Salah Abo-Aba[1,4], Sameera Bafeel[1], Talal A. Zari[1], Sherif Edris[1,3], Ahmed M. Shokry[1,5], Nour O. Gadalla[1,6], Ahmed M. Ramadan[1,5], Ahmed Atef[1], Magdy A. Al-Kordy[1,6], Fotoh M. El-Domyati[1,3], Robert K. Jansen[1,2*]

(Sabir et al., 2014)
PMID: 24718264

# DDBJ pipeline : NGS read alignment

http://p.ddbj.nig.ac.jp/



1) Create new account

2) Login

3) English manual

# DDBJ pipeline 2: Import SRA data



**Selecting Query Files**

1)Click the panel

FTP upload | Private DRA entry | **Import public DRA** | Preprocessing | HTTP upload

Import public FASTQ files from DRA database.

Here is do the section of automatic download of public DRA/ERA/SRA entries.

**Please input DRA/ERA/SRA accession number.** Then the pipeline system import metadata and FASTQ files from DRA database.

Input **DRA/ERA/SRA** Accession Number

SRA100551 | Add my DRA entry

2) Input SRA accession id and click the butt

Accession Number can find here.
DRA Search

Your request. (Here is display only. can not select.)

To select your downloaded entries. See Private DRA entry tab.
When the status makes "done", your requested entry is added in "Private DRA entry" tabs.
When the status makes "failed" or "preparing", please retry it.

**queued** : waiting or during download, **done** : file is ready, **failed** : please retry it, **preparing** : file is not yet in DRA **unchecked** : download is ok, but md5 was not check.

| Status | Submission | Request date |
| --- | --- | --- |
| ○ queued | SRA100551 | 2015-11-25 15:04:35.272 |

3) Import job status

| Status | Submission | Request date |
| --- | --- | --- |
| ✓ done | SRA100551 | 2015-11-25 15:04:35.272 |

4) The pipeline will send
the e-mail notification after job completed.

# DDBJ pipeline 3: Confirm SRA metadata



5) Click the panel

6) Select "SRA100551" dataset

7) Check SRA sample.xml

8) Select paired

9) Go to next

# DDBJ pipeline 4: Specify the alignment tool and generate 9 query sets

# DDBJ pipeline 5: Specify the reference sequence for read alignment analysis



**Specifying Database of Reference Genome**

RESET  BACK  NEXT

○ Major genome sets

○ User original sets

● Download or upload reference

Retrieving a chromosome from DDBJ-DB by using HTTP REST

Input Accession Number (INSD) or (RefseqID)

GU811709

LOAD

PIPELINE    DDBJ-DB

Request

HTTP REST *

Data (fasta)

INTERNET

* Representational State Transfer(REST)

**14) Input GU811709 and Push LOAD button**

Uploading reference from local drive.

FASTA only  参照...  ファイルが選択されていません。  ● UPLOAD

2GB Filesize Limit

☑ >GU811709|GU811709.2 Phoenix dactylifera chloroplast, complete genome.  DELETE

CREATE DATASET

RESET  BACK  NEXT

**15) Select the button**

**Create Genome Dataset**

| | files |
|---|---|
| ☑ | >GU811709|GU811709.2 Phoenix dactylifera chloroplast, complete genome. |

Please input a genomeset description.

**Genome Dataset name**  >GU811709|GU811709.2 Phoenix dactylifera chloroplast, complete genome.

BACK  CREATE GENOMESET

**16) Erase the head accessions**

**17) Select the button**

○ Major genome sets

● User original sets

Genome sets  Phoenix dactylifera chloroplast, complete genome.

☑ >GU811709|GU811709.2 Phoenix dactylifera chloroplast, complete genome.

○ Download or upload reference

**18) Go to next**

RESET  BACK  NEXT

# DDBJ pipeline 6: Set options and run all jobs

## Setting for Reference Genome Mapping

BACK | NEXT

18) Go to next

### bwa

**Set optional parameters of the paired-end analysis**

**Step1) Convert reference sequence**

bwa index  -a is (for small-size reference) ▼ refgenome.fasta

Options usage (click)

**Step2) Map**

bwa aln -t 4 [                    ] refgenome.fasta query1.fastq(.fasta) > out1.sai
bwa aln -t 4 [                    ] refgenome.fasta query2.fastq(.fasta) > out2.sai
bwa sampe [                    ] refgenome.fasta in1.sai in2.sai query1.fastq(.fasta) query2.fastq(.fasta) >
out.sam

**Step3)'uniq': Remove multiple hits on the genome from out.sam.**

Please choose uniq mode.

- ○ Do not remove any read.
- ○ Retain pairs when both reads mapped uniquely or one of reads mapped uniquely,
- ● Retain pairs when both reads mapped uniquely, and Discard other pairs.
- ○ Retain uniquely mapped reads and discard multiply mapped reads.

## Run Confirmation

BACK | RUN

20) Run Jobs!

**Destination of mail**

When the request is completed, the system sends an email to this address.

ekaminum@nig.ac.jp                                    * Required
Result files will be deleted 60 days after submission.

19) You can change
the e-mail address
for job finish notification

**Reference Genome Map [bwa]**

**Query sets**

Query set1

| PairedOrientation | RunAccession | RunAlias | RowLength | Quality Score1 | Quality Score2 |
|---|---|---|---|---|---|
| paired | SRR974754 | AJW-001Run | | | |

Query set2

| PairedOrientation | RunAccession | RunAlias | RowLength | Quality Score1 | Quality Score2 |
|---|---|---|---|---|---|
| paired | SRR974758 | PER-001Run | | | |

Query set3

| PairedOrientation | RunAccession | RunAlias | RowLength | Quality Score1 | Quality Score2 |
|---|---|---|---|---|---|
| paired | SRR974792 | SUK-A001Run | | | |

Query set4

| PairedOrientation | RunAccession | RunAlias | RowLength | Quality Score1 | Quality Score2 |
|---|---|---|---|---|---|

# DDBJ pipeline 7: Confirm job status and outputs



21) Click to job status

22) Click to outputs

23) Download output files

# Visualizing alignment reads using BAM files

ASSIGNMENT[7]
Confirming detected SNPs at 38,157-38,181 positions of "MOS-A cultivar" in the BAM file using "samtools" tview function and save figures of tview screenshot.
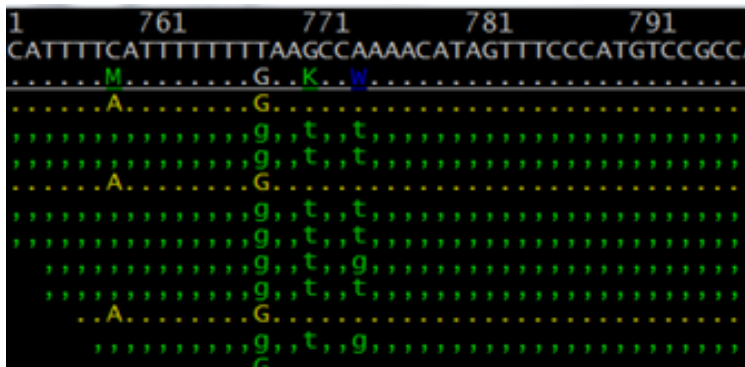
■ Samtools reference
1. http://www.htslib.org/doc/samtools.html
2. https://en.wikipedia.org/wiki/SAMtools

■ Commands in NIG supercomputer
> qlogin
>mkdir datepalm
>cp /home/kaminuma/tmp_data/SRA100551/* ~/datepalm/
>cd datepalm
>less SRR974797_uniqout.sam
>samtools tview SRR974797_out2.bam GU811709.fa

| Cultivar | Position | Reference | Alternate |
|----------|----------|-----------|-----------|
| MOS-A | 38,157 | T | G |
| MOS-A | 38,160 | C | T |
| MOS-A | 38,181 | A | C |

(Table 5. Sabir et al., 2014)

Example : tview screenshot

Reference ➤
Query ➤

# Reference : Alignment file format (SAM/BAM, pileup)

## ＜SAM/BAM format＞

■ Reference
1. https://samtools.github.io/hts-specs/SAMv1.pdf
2. http://genome.sph.umich.edu/wiki/SAM



SAM format (by aligned read)
↓
[compressed]
↓
BAM format

## ＜DDBJ Pipeline download panel＞



1. SRR974797_uniqout.sam
2. SRR974797_out2.bam
3. SRR974797_out2.bam.bai

# Extracting 9 cultivar sequences at psaA/psaB genes by programming

mpileup format

ls -l */SRA100551/pileup/*.pileup

http://samtools.sourceforge.net/mpileup.shtml



Genomic coordinate↑

Reference base↑

Query aligned base↑

Ref. http://www.ebi.ac.uk/ena/data/view/GU811709



Example output : Tab-separated(TSV) file



← 9 cultivar names

1st column: genomic pos.
2nd～10th column:
Aligned base by cultivars

ASSIGNMENT[8]
Extract  9 cultivar genomic sequences from analyzed mpileup files with psaA gene
(genomic position: 40117..42369), and psaB gene (37887..40091) by programming.