

# A manual for large-scale sample collection, preservation, tracking, DNA extraction, and variety identification analysis

Gezahegn Girma, Ismail Rabbi, Adetunji Olanrewaju,  
Oluwafemi Alaba, Peter Kulakow, Tunrayo Alabi,  
Bamikole Ayedun, Tahirou Abdoulaye, Tesfamicheal  
Wossen, Arega Alene, and Victor Manyong





# A manual for large-scale sample collection, preservation, tracking, DNA extraction, and variety identification analysis

Gezahegn Girma, Ismail Rabbi, Adetunji Olanrewaju, Oluwafemi Alaba, Peter Kulakow, Tunrayo Alabi, Bamikole Ayedun, Tahirou Abdoulaye, Tesfamicheal Wossen, Arega Alene, and Victor Manyong

International Institute of Tropical Agriculture, Ibadan

February 2017

Published by the International Institute of Tropical Agriculture (IITA)  
Ibadan, Nigeria. 2017

IITA is a non-profit institution that generates agricultural innovations to meet Africa's most pressing challenges of hunger, malnutrition, poverty, and natural resource degradation. Working with various partners across sub-Saharan Africa, we improve livelihoods, enhance food and nutrition security, increase employment, and preserve natural resource integrity. It is a member of the CGIAR System Organization, a global research partnership for a food secure future.

International address:  
IITA, Grosvenor House,  
125 High Street  
Croydon CR0 9XP, UK

Headquarters:  
PMB 5320, Oyo Road  
Ibadan, Oyo State

ISBN 978-978-8444-83-1

Correct Citation: Girma, G., I. Rabbi, A. Olanrewaju, O. Alaba, P. Kulakow, T. Alabi, B. Ayedun, T. Abdoulaye, T. Wossen, A. Alene and V. Manyong. 2017. A manual for large-scale sample collection, preservation, tracking, DNA extraction, and variety identification analysis. IITA, Ibadan, Nigeria. ISBN 978-978-8444-83-1. 32pp.

Printed in Nigeria by IITA

**Cover photo:** Various stages in the collection of DNA leaf samples, preservation of samples and preparation for DNA sequencing.



# Contents

Acronyms and Abbreviations .....	v
Executive Summary.....	vii
Acknowledgment .....	viii
Introduction.....	1
Establishment of sample tracking system .....	2
Preparation of sample collection kit.....	2
Preparation of self-adhesive stickers of desirable size.....	2
Preparation of duplicate barcode labels .....	4
Booklet and tablet computer for field data entry.....	4
Field sample collection .....	6
Sample collection and preservation .....	6
GPS coordinates and area measurement.....	7
User Guide for Garmin Etrex 10, 20, or 30 GPS devices .....	7
Marking and saving waypoints .....	8
Using the area calculation tool .....	10
Downloading GPS data to your computer .....	11
Troubleshooting the most common problems of GPS devices.....	12
Preparation of samples for DNA extraction .....	14
Implementation of a high throughput DNA extraction method.....	18
Chemicals commonly used for DNA extraction .....	18
Most common recipes required for DNA extraction.....	19
Tris-HCl (1 M, pH 8.0, 100 ml).....	19
EDTA (0.5 M, pH 8.0, 100 ml).....	19
10 x TE-buffer (pH 8.0, 100 ml) .....	19
10 x TBE buffer (1 liter).....	20
25 x TAE buffer (1 liter).....	20
6 x loading dye without xylencyanol (50 ml).....	20
5 M NaCl (100 ml) .....	20
CTAB 2%, medium salt (100 ml).....	20
Decontamination solution (for Ethidium bromide).....	21
Equipment and consumables .....	21
Preparation of DNA extraction buffers.....	21
DNA isolation.....	22
DNA quality and quantity .....	22
Agarose gel electrophoresis .....	23
DNA quantity assessment .....	23
Restriction digestion .....	23
Laboratory Information Management Systems (LIMS).....	25
DNA sample preparation and shipment for genotyping to other labs.....	25

Variety identification analysis.....	27
Genotyping by sequencing.....	27
SNP discovery and quality control of the SNP data .....	27
Establishment of distance threshold to determine identical sets of genotypes .....	28
Cluster analysis .....	29
ADMIXTURE analysis .....	29
Development of a reference library for variety identification .....	29
Matching samples to those in the reference library.....	30
References .....	31
Web resources .....	32

## Tables

1. GPS coordinates and area measurements for varieties of cassava collected from different farmers' fields using GPS devices.....	7
2. The most common problems in Garmin Etrex 10, 20, or 30 devices and recommended solutions.....	13
3. An example of a spreadsheet sample file showing plate name, sample order (numeric), sample coordinate (row, column) and sample name (as on the collection tube). .....	14
4. Chemicals commonly used for DNA extraction with their respective molecular weights.....	18
5. Comparing the nanodrop reading concentration and fluorometer concentration. ....	23
6. Hind III digestion master mix. ....	24

## Figures

1. Example of barcode labels prepared in duplicates for two different samples.....	4
2. Data collection sheet of a booklet with the sample-related information and other metadata: an example from the CMS project.....	5
3. Recommended size and age of cassava leaf tissue, silica gel dried leaf, and sample identification.....	7
4. Garmin Etrex 20 GPS device and main uses of the buttons. ....	8
5. A standard plate map.....	14
6. Samples arrangement on carton made in-house. ....	17
7. DNA quality and quantity assessment including test gel electrophoresis of DNA samples (top) and restriction enzyme digested for randomly selected DNA samples .....	25
8. Samples arranged in a stack (with barcode labels pasted on each plate) in addition to plate ID information. ....	26
9. Pairwise genetic distance (IBS) calculated using SNP markers. ....	28
10. Example of cluster analysis showing genetically identical varieties. ....	29

# Acronyms and Abbreviations

$\beta$	Beta
CMS	Cassava monitoring survey
CGIAR	Consultative Group for International Agricultural Research
CTAB	Cetyltrimethyl ammonium bromide
$^{\circ}\text{C}$	Degrees celsius
DNA	Deoxyribonucleic acid
EDTA	Ethylene diaminetetra acetic acid
g	Gram
GBS	Genotyping by sequencing
GDF	Genomic diversity facility
GPS	Global positioning system
g/mol	Gram per mole
HCl	Hydrogen chloride
IBS	Identity by state
ID	Identification
IITA	International Institute of Tropical Agriculture
kPa	Kilo Pascal
LIMS	Laboratory information management system
M	Molarity
MAF	Minor allele frequency
MB	Megabyte
$\mu\text{l}$	Microliters
ml	Milliliters
mM	Millimolar
$\text{m}^2$	Meter square
NB	Note well
NaCl	Sodium chloride
NE	New England
NGS	Next generation sequencing
ng	Nanogram
NaOH	Sodium hydroxide
%	Percent
PCR	Polymerase chain reaction
KAc	Potassium acetate
PH	Potential hydrogen.
PVP	Polyvinylpyrrolidone
QC	Quality control

RE	Restriction enzyme
rpm	Revolution per minute
RNase	Ribonuclease
SDS	Sodium dodecyl sulfate
TAE	Tris-acetate-EDTA (Ethylenediamine tetra acetic acid)
TBE	Tris-borate-EDTA
TE	Tris-EDTA
UV	Ultraviolet
w/v	Weight per volume

# Executive Summary

Several alternative options have been used for varietal identification. However most of the traditional methods have inherent uncertainty levels and estimates often have wide confidence intervals. In an attempt to circumvent traditional survey-based measurement errors in varietal identification, DNA-based varietal identification has been implemented in the Cassava Monitoring Survey (CMS) of Nigeria — a large adoption study involving 2500 cassava farming households. The DNA fingerprinting technique offers a reliable method to accurately identify varieties grown by farmers and increases accuracy and credibility in the interpretation of adoption rates and associated economic and policy analyses. Unlike phenotype-based methods, DNA is not affected by environmental conditions or plant growth stage and is more abundant than morphological descriptors. However, undertaking a credible DNA-based varietal identification is not a trivial matter because of the logistical challenges involving sample collection and tracking by a large team of field enumerators. This manual presents the detailed steps required for undertaking reliable DNA-fingerprinting-based identification of cassava varieties. In particular, the manual gives detailed information on the establishment of a sample tracking system, preparation of a readily available and cheap sample collection kit, field sample collection methodology, preparation of samples for DNA isolation, and development of a pipeline for variety identification analysis. This manual is part of the outputs of the CMS project funded by the CGIAR Research Program on Roots, Tubers and Bananas (RTB), the Bill & Melinda Gates Foundation, and the International Institute of Tropical Agriculture (IITA).

# Acknowledgment

The CGIAR Research Program on Roots, Tubers and Bananas (RTB) and the Bill & Melinda Gates Foundation (Gates Foundation) funded the Cassava Monitoring Survey (CMS) project. The authors acknowledge Dr Alene Arega, Dr Shiferaw Feleke, Dr Abass Adebayo and Mr Henry Musa Phaka for their contributions to the conceptualization and design of the CMS project. They are grateful to Elvis Fraser for his guidance in the design of this study. They thank all the participants at the Nigeria Cassava Monitoring survey convening organized by Bill & Melinda Gates Foundation in Dar es Salaam, Tanzania from 15–21 March 2015. They thank all other staff of the IITA Socioeconomics Unit and Bioscience Center for their contributions to the success of the project and development of this manual. The contribution of the National Roots Crops Research Institute (NRCRI), Umudike, Nigeria to the success of the CMS Project is also acknowledged. The contributions of enumerators, extension agents and most importantly Cassava farmers in Nigeria is gratefully acknowledged.

# Introduction

DNA fingerprinting offers a reliable method to accurately identify biological samples. It is applied in all areas of the biological sciences in general and crop improvement in particular including germplasm characterization, cultivar identification, and genetic diversity studies (Rabbi et al., 2015). However, the efficiency of DNA fingerprinting depends mainly on the quality of the tissues (i.e., well-preserved leaf samples) for the recovery of high quality DNA and a good sample tracking system to ensure the chain of custody of samples for identification from the source (e.g., experimental field, farmer's field, and screen house) to the laboratory.

This particular manual was developed based on the experience from a project entitled Cassava Monitoring Survey (CMS), a study designed to assess the drivers of adoption of improved cassava cultivars in Nigeria. The DNA fingerprinting of cassava cultivars collected from 2500 households was one of the main components of the project. Genotyping-by-sequencing (GBS) procedure was applied for the identification of cultivars by matching them to a "library" consisting of known improved cultivars, landraces, and the genebank collection previously genotyped with a similar procedure (Rabbi et al. 2015). High quality DNA is the main requirement for GBS. The main activities in the DNA fingerprinting component of the project included the following: sample kit preparation, establishment of tracking system, field tissue collection, Global Positioning System (GPS) coordinates and area measurements, high throughput DNA extraction, genotyping, and variety identification.

The experience from this study will hopefully help other similar studies involving large-scale DNA-based adoption studies. The manual presents detailed steps for the establishment of a sample tracking system, preparation of a sample collection kit, field sample collection, preparation of samples for DNA isolation, and development of a pipeline for variety identification analysis. In addition protocols are presented for large sample size DNA extraction including preparation of all the required extraction buffers, and the preparation of the most common recipes as well as important considerations for DNA sample storage and shipping. This manual is part of the outputs of the CMS project funded by the CGIAR Research Program on Roots, Tubers and Bananas (RTB), the Gates Foundation, and IITA.

# Establishment of sample tracking system

A standard tracking system is important particularly when we are dealing with a large sample size collected by dozens of field enumerators to reduce any possible introduction of human errors of sample mismatch and mix up. The use of multiple layers of tracking system will improve the accuracy and reliability of the study in general. Here we introduce a simple and cheap sample tracking system that was found to be very useful and effective for our previous study. The items required and important steps followed for establishment of the tracking system are detailed as follows.

## **Preparation of sample collection kit**

A sample collection kit consisting of plastic tubes, silica gel, adhesive labels, barcode labels, and printed and bound booklets needs to be prepared before the field visits. Silica gel is a granular, vitreous, porous form of silicon dioxide made synthetically from sodium silicate. As a desiccant, it has an average pore size of 2.4 nanometers and a strong affinity for water molecules.

The silica gel weight requirement varies depending on the tissue type. For example, for collection of two newly expanded young leaves of cassava, approximately 20 g of silica gel in a 50-ml plastic vial was found to be sufficient. It is preferable to use colored silica gel desiccant that changes when it is hydrated and thereby unable to function as expected. Likewise, the amount of silica gel required can be optimized for different plant tissues. Prior to field visits it is also advisable to paste pre-prepared labels (self-adhesive stickers) on the vials with enough space for information to be captured.



### **Precautions:**

Silica gel is a non-hazardous chemical but it may cause the hands to become dry. Please always make sure gloves are worn to reduce direct hand contact leading to such a situation.

### *Preparation of self-adhesive stickers of desirable size*

We advise researchers to use printed self-adhesive labels for recording information on the sample jars. For the CMS project we adopted the Phenix self-adhesive permanent labels ([www.phenixresearch.com](http://www.phenixresearch.com)). Following is a step-by-step approach for formatting such labels in Microsoft word for laser printers.

1. Select Labels from the Tools menu (for XP users select Letters and Mailings from the Tools menu then select Envelopes and Labels).
2. Click on the Labels tab and choose Options.
3. Click on New Label and type the name of the label you want to format (e.g., LBS-1000L or another name) in the space for Label name.
4. Input the following parameters as shown below in appropriate boxes and press OK.

**New Custom laser**

Preview

The diagram shows a label sheet with a grid of labels. The dimensions are labeled as follows:
 

- Top margin:** The space between the top edge of the sheet and the top edge of the first row of labels.
- Side margins:** The space between the left and right edges of the sheet and the left and right edges of the first column of labels.
- Horizontal pitch:** The distance between the center of one label and the center of the next label in the same row.
- Vertical pitch:** The distance between the top edge of one label and the top edge of the label directly below it.
- Width:** The width of a single label.
- Height:** The height of a single label.
- Number down:** The total number of labels in a column.
- Number across:** The total number of labels in a row.

Label name:

Top margin:  Label height:

Side margin:  Label width:

Vertical pitch:  Number across:

Horizontal pitch:  Number down:

Page size:

5. Go back to Options and highlight the selection you need (or others) and press OK.
  6. For single labels or a full sheet with the same information, type it in this dialog box or click on New Document to type information directly on each label. Print. (N.B. Centering the text is recommended.)
- LBS-1000L 2 5/8" × 1" labels have a format identical to Avery 5160, 5260, and 8160 labels.

### *Preparation of duplicate barcode labels*

Barcode labels provide unique identifiers for each sample and serves as a backup sample tracking layer when one copy is pasted on the sample vial and the other copy next to the recorded sample in the survey booklet (see section 2.2 below). The list below is a step-by-step guide for the preparation of barcode labels.

1. Download IDAutomation barcode software that provides barcode fonts from <http://www.idautomation.com/barcode-fonts/code-39/download.html>.
2. Open Microsoft Excel and enter a series of numbers corresponding to sample size in duplicate. Different formats can be used, e.g., CMS-10001, \*10001\*, CMS10001.
3. Create rows and columns according to the desirable size on Microsoft Word. Row by column of 21 and 8 respectively is recommended for 84 different samples in duplicates.
4. Change the font style to IDAutomationHC39M from font menu and adjust the font size so the bars and identification under the bars can be clearly seen, as shown below (Fig.1).
5. Adjust the printer settings to labels and then print.



**Figure 1: Example of barcode labels prepared in duplicates for two different samples.**

### **Booklet and tablet computer for field data entry**

A hard copy booklet detailing correct methods on sampling and data collection sheets for each household is important for recording sample-related information and ensuring that information is captured at multiple levels. For example, in the CMS project, similar information on stickers that includes the IDs for region, enumeration area, and household, as well as variety ID and name was also captured on booklets (Fig. 2). Additional information was collected including the household head's name, GPS co-ordinates of where the questionnaire was used and on the farm where the sample was collected, field and plot IDs, plot size, and cropping pattern (intercropped/ monocropped). There was also a space for pasting the barcode label identical to the one put on the sample collection tube for all the different cultivars. The disadvantages of using a booklet are the time and personnel requirements and the possible introduction

of errors during data entry. However, it can still be used as a backup and another important layer in a sample tracking system.

Remarks:

- In the CMS project or for any other related study addressing variety identification for adoption study, the “VARIETY ID” and “VARIETY NAME” should be completed before you proceed to the field. Once in the field, the farmer will be asked to point out the varieties he/she named and then the remaining sample information can be recorded.
- The barcode number on the sample collection tube should be typed in the survey instrument next to the sample identifier.
- The GPS co-ordinates should be captured for each field visited.

<b>Household information</b>	This corresponds with the information provided by farmer		Here you fill the GPS co-ordinates where questionnaire interview is taking place
Region	SOUTH WEST	Latitude	7.3963889 N
State	OYO	Longitude	3.9166667 E
Local government area	AKINYELE	Sample collector's name	Mr BANKOLE OYEWALE
Enumeration area ID	111011	Date of sample collection	14.06.2015
Household ID	1110112	Supervisor's name	Dr OLANIYI ABIODUN

### Sample information table

Var. ID	Variety name	Field	Plot	Area (m <sup>2</sup> )	GPS		Planting	BARCODE
V1	Oko-lyawo	1	1	956.5	Latitude	7.3963889 N	Monocropped	*10001*
					Longitude:	3.916944 E 4E		
V2	Agric	2	1	245.2	Latitude:	7.3963889 N	Monocropped	*10002*
					Longitude:	3.9169456 EE		
								

Barcode stickers are provided in duplicates. One copy should be placed here and the second one on the corresponding sample tube.

**Figure 2: Data collection sheet of a booklet with the sample-related information and other metadata: an example from the CMS project.**

The use of a tablet computer for the entry of sample-related information has an advantage, as the data can easily be transferable unlike the use of a booklet. However, there could be a risk of system crash that can lead to a complete loss of data. Hence, a regular backup to external disks is needed. The prices of tablet computers are reducing but they are still costly for research in developing countries.

## **Field sample collection**

Field sample collection includes collecting and recording leaf tissue and recording sample-associated information. Proper sampling and conservation of plant tissues and the capture of sample-related information are essential to assure the quality of DNA and associated information. It is the critical step for the success of genotyping work. In some cases sampling involves remote field collection that it is usually very difficult to resample. We recommend collecting at least twice the amount of leaf tissue needed for DNA extraction. It is therefore very important to give maximum emphasis during field sample collection and DNA extraction to reduce the risk of sample loss. The backup leaf tissue can be used in case of sample loss during the extraction process.

## **Sample collection and preservation**

Ensuring tissue is collected of appropriate size and age is very important for a large quantity of high quality nucleic acids to be extracted. In addition it is also important to make sure that the sample to be collected represents the intended individual/variety. Hence, intensive training of the personnel involved in sample collection is critical prior to field visits. Below are the main steps to be followed once you are ready for sample collection.

1. For cassava, collect two newly expanded apical leaf tissues of approximately 6 cm from a single stem (Fig. 3) representing the exact intended sample and place them inside the tube containing silica gel. Likewise, the requirements for tissue age and size can be optimized for other plant species depending on the quality and quantity of DNA yield.
2. Enter all the required information on the label pasted on the surface of the sample tube and also place duplicate barcode labels distinct for all the samples, one on the tube and the other on the booklet (Fig. 3). In parallel, enter the barcode label codes on the Surveybe (a data collection software) next to the variety identification (variety ID and name).
3. Turn the sample collection tube upside down to make sure that all the leaf tissues are in contact with the silica gel for complete desiccation.
4. At the same time, write down the information associated with the samples on a booklet and Surveybe questionnaire. Likewise the same information can be captured electronically using tablets if available. We strongly recommend recording the barcode number in the Surveybe questionnaire for the corresponding variety.



Figure 3: Recommended size and age of cassava leaf tissue, silica gel dried leaf, and sample identification

### GPS coordinates and area measurement

GPS is a system to estimate location on earth by using signals from a set of orbiting satellites. GPS field coordinates and plot size need to be captured while leaf tissues of individual accessions representing different cultivars (in all farms visited) are being collected (Table 1).

**Table 1: GPS coordinates and area measurements for varieties of cassava collected from different farmers' fields using GPS devices.**

Variety ID	Variety name	Field ID	Plot ID	Area (m <sup>2</sup> )	GPS		
V1	OKO IYAWO	1	1	1,124.3	Longitude	004	52.579'
					Latitude	07	32.453'
V2	IDILERU	2	1	6,102.9	Longitude	004	52.583'
					Latitude	07	32.456'

### User Guide for Garmin Etrex 10, 20, or 30 GPS devices

General description of Garmin Etrex 20 and main uses of the buttons (Fig. 4)

1. Move the **Thumb Stick** up, down, left, right, to highlight menu selections or move around the map.
2. Press the center of the **Thumb Stick** item to select the highlighted item.
3. Press **back** to move back one step in menu structure.

4. Press **Menu** to display a list of functions/options for the current page.
5. Press **Menu** twice to access the main menu from any page.
6. Press **Up** and **Down** arrow to zoom in and out on the map.

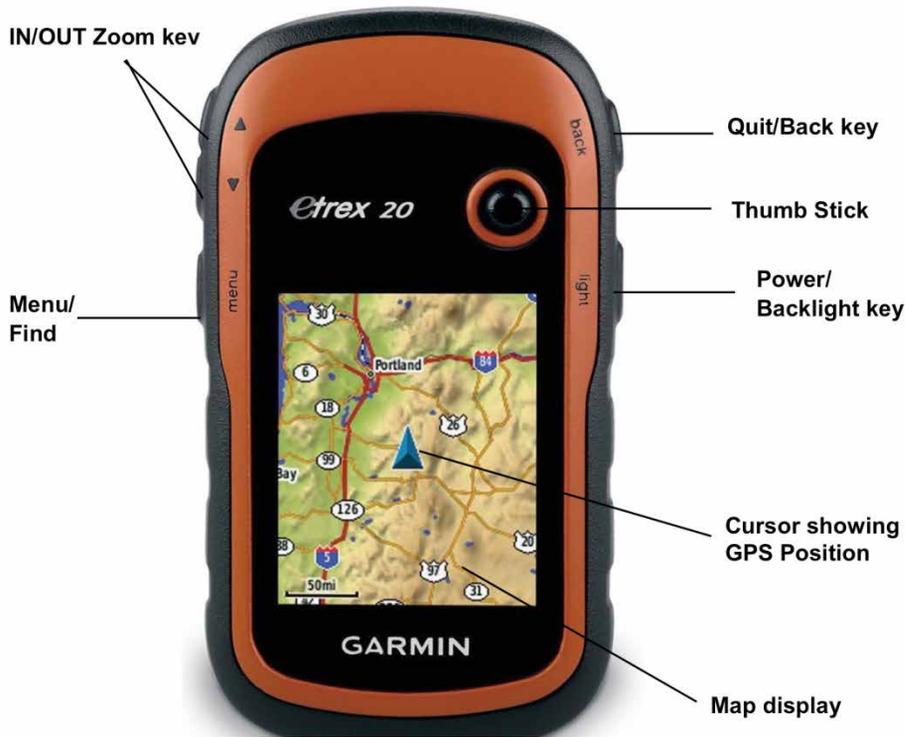


Figure 4: Garmin Etrex 20 GPS device and main uses of the buttons.

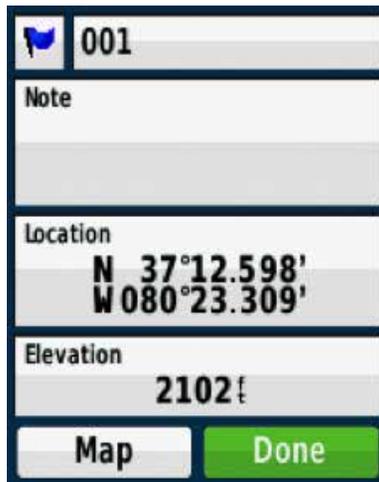
### *Marking and saving waypoints*

A **waypoint**, otherwise known as a landmark, is a reference point in physical space used for purposes of navigation. Examples of waypoints in the CMS are the location of the village, cassava farm, and farmer's household. Once a waypoint is established and saved you can easily navigate back to it using the GPS. This can be helpful in locating research plots or a cassava farm, identifying and relocating potential pollution sources or identifying and marking specific places within an area that have been affected by blight or disease (where you may want to go again at a later date). It is also possible to upload waypoints to a desktop computer (using DNR Garmin or GPS Utility or a similar software program).

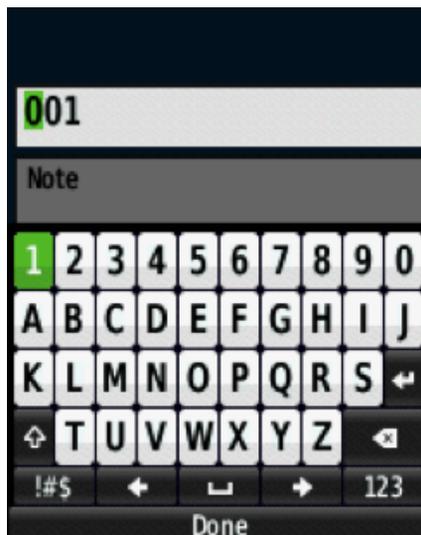
To mark and save your current location as a waypoint:

1. Walk to the point where you want to obtain a waypoint (a farmer's homestead, a cassava farm, a village center, etc.).

2. Press down and hold the **click-stick** until the waypoint page appears (see the picture below). You can also go to **Main Menu** and select **Mark Waypoint**.



3. This GPS automatically assigns three digit numbers to waypoints (in this example, it assigned **001** as the waypoint name). You can customize the name of the waypoint. To change the name of the waypoint, use the **click-stick** to highlight the waypoint name field (in the above picture, the name field is **001**) and click **Straight down**.
4. Type the new name for the waypoint, using the **click-stick** to select and enter the characters from the on-screen keyboard and select **Done** (in small print, at the bottom of the screen) as shown below.



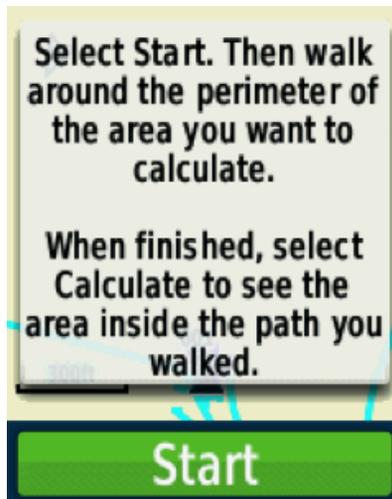
5. Use the **click-stick** to scroll down to the **Done** button (refer to the above picture). When you select this it will save the waypoint and take you back to the main menu.

*Using the area calculation tool*

1. Generally you should clear the current active track log just before you begin laying tracks. Go to: Main Menu > Track Manager > Current Track > Clear Current Track. [If you wish to save the previous track log, first select SAVE then clear the log.]
2. Go to Main Menu >Area Calculation (see below).

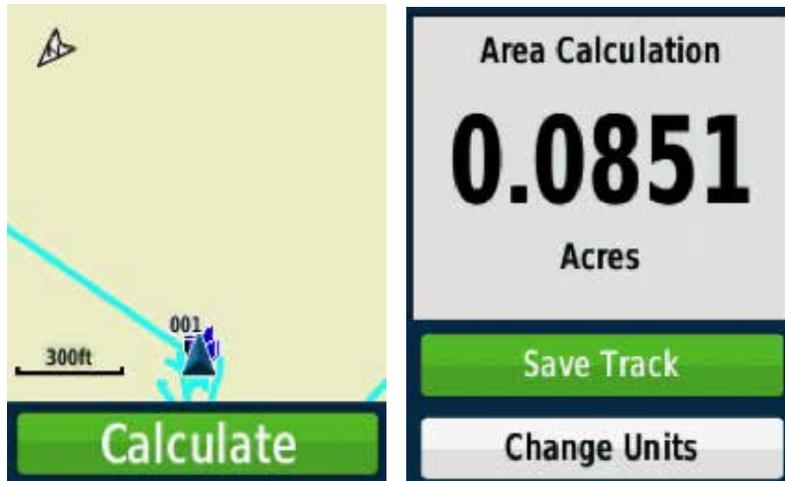


3. The Area Calculation page will have a Start option at the bottom (see the picture below). Once you are at your starting point, click Start using the click-stick.



4. Begin walking around the perimeter of the area that you want to calculate. The screen displays your progress. Zoom in or out as appropriate to view your tracks. Don't despair if heavy tree cover causes you to occasionally lose contact with the satellites as you track. The GPS will "connect the dots" and link your recorded track points in an attempt to estimate the enclosed area. View the saved track screen to decide whether or not the integrity of the track was maintained.

- Once you return to your starting point, click on Calculate (see below). If you aren't at your starting point, your receiver will automatically complete the loop with a straight shot from your current position to your starting point. The enclosed area value will be displayed (see the picture below). To change units, highlight and click on the Change Units option to bring up a selectable list of choices (square feet, square yards, square meters, hectares, square miles, etc)



- If this area calculation is something you'll need to refer to later, select **Save Track**. A page will open that will allow you to re-name the track if necessary.
- To view all of your saved tracks, go to: **Main Menu > Track Manager > Archived Tracks**.

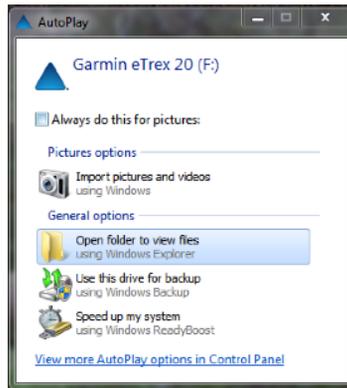
### *Downloading GPS data to your computer*

The Garmin Etrex 10, 20, or 30 are USB compatible devices and each of them comes with a USB cable.

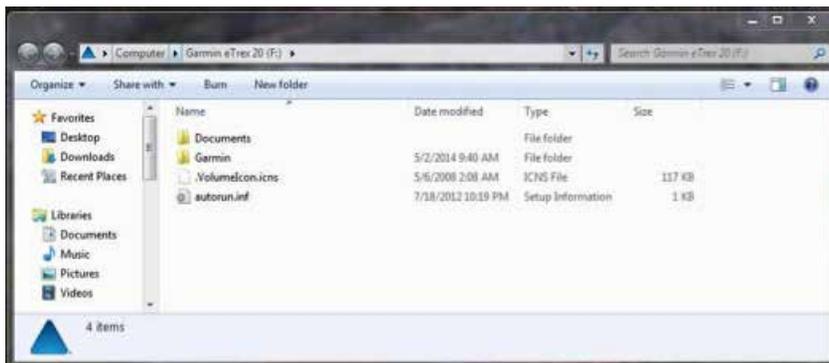
- Connect the Gamin Etrex 20 to your computer using the supplied USB cord. Once the unit is connected to your computer it should have a USB icon on the screen as shown below.



- When the GPS is connected and your Windows computer recognizes it, you should see the following dialog box.



- Select "Open folder to view files". You should see something like this.



- Open the Garmin folder such as this following path I:\Garmin\GPX\
- Copy Tracks or Waypoints files with extensions \*.gpx to your computer
- Disconnect your GPX unit by right clicking the "Safely Remove Hardware and Eject Media" icon on the task bar. Select "Eject Garmin Etrex 20".

### *Troubleshooting the most common problems of GPS devices*

You can also troubleshoot the most common problems in GPS devices by following very simple steps (Table 2). If the problems are beyond what is listed here further information can be obtained through browsing web resources or asking GIS experts.

**Table 2: The most common problems in Garmin Etrex 10, 20, or 30 devices and recommended solutions.**

Problem	Solution
The device does not respond. How do I reset the device?	Remove the batteries. Reinstall the batteries. Note: This does not erase any of your data or settings.
I want to reset all settings back to the factory defaults.	Select <b>setup &gt; Reset&gt; Reset All settings</b>
My device does not acquire satellite signals.	Take your device out of buildings and away from tall buildings or trees. Turn on the device. Remain stationary for several minutes.
How do I know my GPS is in USB mass storage mode?	On the device, an image of the device connected to a computer appears. On your computer, you should see a new removable disk drive in My computers on Windows computers.

# Preparation of samples for DNA extraction

Once samples are received in the laboratory the first step should be the arrangement of the sample collection tubes in a set of 96 well plates (Figs 5 and 6). In each plate it is advisable to include a blank sample as a control and a replicate of one random sample for quality control (QC). In the CMS project a blank sample was kept in a well in such a way that it matched with the plate number for easy tracking. For example, the 20<sup>th</sup> plate labeled as CMS-Plate-020 had a blank sample in the 20<sup>th</sup> well. The labeling of each sample was also made following the standard plate map well order (Fig. 5, Table 3).

Sample collection date and method: \_\_\_\_\_

Plate Name: \_\_\_\_\_

Sample information: \_\_\_\_\_

	1	2	3	4	5	6	7	8	9	10	11	12
A												
B												
C												
D												
E												
F												
G												
H												

Figure 5: A standard plate map.

Table 3: An example of a spreadsheet sample file showing plate name, sample order (numeric), sample coordinate (row, column) and sample name (as on the collection tube).

Plate name	Well number	Well ID	Plastic vial (sample collection tube) ID
Plate 001	1	A01	BLANK
Plate 001	2	B01	001-2
Plate 001	3	C01	001-3
Plate 001	4	D01	001-4
Plate 001	5	E01	001-5
Plate 001	6	F01	001-6
Plate 001	7	G01	001-7
Plate 001	8	H01	001-8

Plate name	Well number	Well ID	Plastic vial (sample collection tube) ID
Plate 001	9	A02	001-9
Plate 001	10	B02	001-10
Plate 001	11	C02	001-11
Plate 001	12	D02	001-12
Plate 001	13	E02	001-13
Plate 001	14	F02	001-14
Plate 001	15	G02	001-15
Plate 001	16	H02	001-16
Plate 001	17	A03	001-17
Plate 001	18	B03	001-18
Plate 001	19	C03	001-19
Plate 001	20	D03	001-20
Plate 001	21	E03	001-21
Plate 001	22	F03	001-22
Plate 001	23	G03	001-23
Plate 001	24	H03	001-24
Plate 001	25	A04	001-25
Plate 001	26	B04	001-26
Plate 001	27	C04	001-26_QC
Plate 001	28	D04	001-28
Plate 001	29	E04	001-29
Plate 001	30	F04	001-30
Plate 001	31	G04	001-31
Plate 001	32	H04	001-32
Plate 001	33	A05	001-33
Plate 001	34	B05	001-34
Plate 001	35	C05	001-35
Plate 001	36	D05	001-36
Plate 001	37	E05	001-37
Plate 001	38	F05	001-38
Plate 001	39	G05	001-39
Plate 001	40	H05	001-40
Plate 001	41	A06	001-41
Plate 001	42	B06	001-42
Plate 001	43	C06	001-43
Plate 001	44	D06	001-44
Plate 001	45	E06	001-45
Plate 001	46	F06	001-46
Plate 001	47	G06	001-47
Plate 001	48	H06	001-48
Plate 001	49	A07	001-49
Plate 001	50	B07	001-50
Plate 001	51	C07	001-51
Plate 001	52	D07	001-52
Plate 001	53	E07	001-53
Plate 001	54	F07	001-54

Plate name	Well number	Well ID	Plastic vial (sample collection tube) ID
Plate 001	55	G07	001-55
Plate 001	56	H07	001-56
Plate 001	57	A08	001-57
Plate 001	58	B08	001-58
Plate 001	59	C08	001-59
Plate 001	60	D08	001-60
Plate 001	61	E08	001-61
Plate 001	62	F08	001-62
Plate 001	63	G08	001-63
Plate 001	64	H08	001-64
Plate 001	65	A09	001-65
Plate 001	66	B09	001-66
Plate 001	67	C09	001-67
Plate 001	68	D09	001-68
Plate 001	69	E09	001-69
Plate 001	70	F09	001-70
Plate 001	71	G09	001-71
Plate 001	72	H09	001-72
Plate 001	73	A10	001-73
Plate 001	74	B10	001-74
Plate 001	75	C10	001-75
Plate 001	76	D10	001-76
Plate 001	77	E10	001-77
Plate 001	78	F10	001-78
Plate 001	79	G10	001-79
Plate 001	80	H10	001-80
Plate 001	81	A11	001-81
Plate 001	82	B11	001-82
Plate 001	83	C11	001-83
Plate 001	84	D11	001-84
Plate 001	85	E11	001-85
Plate 001	86	F11	001-86
Plate 001	87	G11	001-87
Plate 001	88	H11	001-88
Plate 001	89	A12	001-89
Plate 001	90	B12	001-90
Plate 001	91	C12	001-91
Plate 001	92	D12	001-92
Plate 001	93	E12	001-93
Plate 001	94	F12	001-94
Plate 001	95	G12	001-95
Plate 001	96	H12	001-96

QC = quality control sample replicated from a randomly selected sample (B04 replicated in CO4 in the above example).

The next step is to capture information pasted on the sample collection tube on a hard copy plate map (Fig. 5) or a notepad using DNAplateApp software. Information from stickers on each sample collection tube can be captured on the notepad and manually on a hard copy plate map, using a barcode reader for the barcode label. A DNAplateApp software, freely accessible from the Poland Lab for Wheat Genetics at Kansas State University (<http://wheatgenetics.org/research/12-research/software/21-dna-plate-app>) with a detailed information guide, can be used to capture the barcode label using a barcode reader/scanner.

Once all information is captured the step to be followed is to transfer the optimum sized tissues (e.g., approximately 200 mg silica gel dried leaf tissue for cassava) into 1.2 ml sized extraction tubes arranged on a rack. Make sure two stirring magnetic balls are inserted inside each extraction tube prior to sample transfer. The important considerations during sample transfer are avoiding mix up and making sure equal volumes of samples are added. It is therefore advisable to take each strip of eight extraction tubes on a separate rack while transferring samples and compare the sample size across different tubes for consistency.

Finally transfer the plastic vials with remaining samples in a similar order on a permanent structure of 30cm x 45cm size and assign them distinct plate numbers (Fig. 6).

NB: Similar structures can be made from different items such as plywood.



**Figure 6: Samples arrangement on carton made in-house.**

# Implementation of a high throughput DNA extraction method

We have implemented a modified high throughput DNA extraction protocol (Rabbi et al. 2014) that enables extraction of up to 10 plates (960 samples) at a time. It is always advisable to list the required chemicals, prepare reagents with the appropriate concentration, and make sure about the availability of the basic laboratory equipment and consumables before commencing on the isolation of nucleic acid. The most common recipes important for DNA extraction have been also listed (section 4.2).

## Chemicals commonly used for DNA extraction

Below are lists of chemicals commonly used for DNA extraction, purification, and quality and quantity assessments with their respective molecular weights (Table 4).

### Precautions:

- Always make sure you are up to date on the safety precautions by reading the appropriate Material Safety Data Sheets first before starting to work with any chemical.
- It is also important to always check and confirm the exact name, molecular weight, and expiry date of any chemical before use so as to make sure that the chemical used is the one intended and still valid.

**Table 4: Chemicals commonly used for DNA extraction with their respective molecular weights.**

Chemical	Molecular weight [g/mol]
Agarose	630,54
Ammonium acetate	77,08
Boric acid	61,83
Bromphenol blue	691,90
Chloroform	119,38
CTAB	364,46
EDTA	372,20
Ethanol	46,07
Ethidium bromide	394,29
Glycerol	92,09
HCl	36,46

Hypophosphorous acid	66,00
Isopropanol	60,10
Isoamyl alcohol	88,15
NaCl	58,44
NaOH	40,01
PVP	2,50
KAc	266,12
Sodium nitrite	68,99
Tris	121,10
$\beta$ -Mercaptoethanol	78,13
Sodium dodecyl sulfate	288,37
Sodium nitrite	69,00

## Most common recipes required for DNA extraction

### Remark: v-

Please always work in a hygienic workspace. Make sure of this by cleaning the bench, micropipettes, and other plastic ware with 70% ethanol.

### *Tris-HCl (1 M, pH 8.0, 100ml)*

- Dissolve 12.1 g Tris in approx. 70ml demineralized water
- Adjust the pH to 8.0 with concentrated HCl
- Mix well and add demineralized water up to 100 ml
- Autoclave at 121 °C
- Store at room temperature

### *EDTA (0.5 M, pH 8.0, 100ml)*

- Dissolve 18.6 g EDTA in approx. 70 ml of demineralized water
- Adjust the pH to 8.0 with NaOH (e.g., 1M solution) EDTA dissolves only after the adjustment of the pH
- Mix well and add demineralized water up to 100 ml
- Autoclave at 121 °C
- Store at room temperature

### *10 x TE-buffer (pH 8.0, 100ml)*

- Mix 1 ml Tris-HCl 1M and 200  $\mu$ l EDTA 0.5M
- Mix well and add demineralized water up to 100 ml
- Autoclave at 121 °C
- Store at 4 °C

### *10 x TBE buffer (1 liter)*

- Dissolve 108 g Tris and 55 g Boric acid in approximately 700 ml of sterile demineralized water
- Add 40 ml 0.5M Na-EDTA
- Adjust pH to 8.57 with 5M NaOH
- Mix well and add demineralized water up to 1 liter
- Store at room temperature

### *25 x TAE buffer (1 liter)*

- Dissolve 121 g Tris in approx. 700ml of sterile demineralized water
- Carefully add 57.1ml glacial acetic acid in a fume hood
- Add 50 ml 0.5M Na-EDTA
- Mix well and add demineralized water up to 1 liter
- Store at 4 °C

### *6 x loading dye without xylencyanol (50 ml)*

- 0,125 g Bromophenol blue
- 1 ml Tris-HCl
- Dissolve in 34 ml of demineralized water
- Add 15 ml Glycerol (= 1.26 \* 15 = 18.9 g)
- Make aliquots in sterile eppendorf caps
- Store at 4 °C

### *5 M NaCl (100 ml)*

- Dissolve 29.22 g NaCl in 100 ml of demineralized water
- Autoclave at 121 °C
- Store at room temperature

### *CTAB 2%, medium salt (100 ml)*

- Add in a 100 ml flask
- 10 ml, 1M Tris-HCl
- 28 ml, 5M NaCl
- 4 ml, 0.5M EDTA
- 2 g CTAB
- 1 g PVP
- Store at 4 °C

Add **sterile** demineralized water up to 100 ml

(It is possible to warm the solution up to 50 °C to speed up the process.)



**DO NOT AUTOCLAVE THIS SOLUTION!!!**

### *Decontamination solution (for Ethidium bromide)*

- Weigh in 0.7 g Sodium nitrite
- Add 3.3 ml Hypophosphorous acid (50%)
- Add demineralized water up to 50 ml

### **Equipment and consumables**

- Weighing balance
- PH meter
- Genogrinder
- Fume hood
- Ultra-centrifuge
- Freezer (-20 °C)
- Ice machine
- Incubator (37 °C)
- Water bath @ 65 °C
- Spectrophotometer
- Electrophoresis apparatus
- Gel documentation,
- Extraction tubes (1.2ml)
- Multiple channel pipette
- Single micropipettes
- Pipette tips
- Eppendorf
- 96-well PCR plates,
- Extra 96-well holding racks
- Pair of stirring balls per extraction tubes
- Magnets
- Hot plate stirrer
- Stir bar
- Paper towels

### **Preparation of DNA extraction buffers**

The DNA extraction protocol described below is a modification from Dellaporta et al. (1983) for DNA isolation from large samples of up to 10 plates per day.

- Prior to the extraction of high molecular weight DNA, prepare a large volume (1 liter) of extraction buffers consisting of
  - 1% (w/v) PVP (i.e., 10 g),
  - 100mM Tris-HCl (100 ml of 1 M Tris-HCl, pH 8.0),
  - 50mM EDTA(100 ml of 0.5 M EDTA, pH 8.0),
  - 500mM NaCl (100 ml of 5 M NaCl).

- Stir the mixture using a stir bar; meanwhile adjust the pH to 8.0 using HCl.
- Prepare the mixture up to 1000 mL with deionized water and autoclave for 15 min at 121 °C under a pressure of 100kPa.
- The buffer could be stored for up to six months depending on how it is being used.
- Prepare stock solutions of the different working buffer components such as 20% SDS, 0.75% β-mercaptoethanol following standard procedures.
- Prepare a total volume of 500 ml working buffer from the stock extraction buffer prior to the DNA extraction, with addition of 33 ml of 20% (w/v) Sodium dodecyl sulfate and 0.75% β-mercaptoethanol (375 μl) just before use.

### **Precautions:-**

Working buffer should be prepared in a fume hood to avoid the risk of inhaling dangerous chemicals; the standard laboratory waste-disposal should be strictly adhered to when extraction mixtures and used plastic ware are discarded.

### **DNA isolation**

- Balance two boxes each consisting of 96 samples (including 1 blank) containing extraction tubes on a weighing balance and proceed with grinding to fine powder for 1 min at 1500 rpm using Geno|Grinder SPEX SamplePrep 2010.
- NB: Look at each sample tube to make sure that all the samples are well ground.
- Add 400 μl of extraction buffer (containing freshly prepared β-mercapto- ethanol) to the tissue powder in each extraction tube.
- Incubate the mixture at 65 °C for 10 min with gentle mixing by inversion.
- Remove the tubes from the incubator and immediately add 200 μl of ice-cold 5M Potassium acetate and incubate on ice for 20 min.
- Add 350 μl of Chloroform isoamyl alcohol (24:1) into the side of tubes and centrifuge at 4000 g for 10 min.
- Decant the supernatant into new extraction tubes already containing 400 μl of 100% isopropanol, then centrifuge at 4000 g for 10 min to precipitate the DNA as a pellet.
- Discard isopropanol and wash DNA pellets by adding 300 μl of 70% Ethanol, flapped and then centrifuge at 3500 g for 10 min before discarding the ethanol.
- Air-dry the DNA pellets for a few minutes and suspend in 100 μl of autoclaved low-salt TE buffer.
- Finally, add 10 μl of DNase-free RNaseA to the DNA and incubate at 37 °C for 2 hr.
- Store the DNA at –20 °C for future use.

### **DNA quality and quantity**

For construction of a successful genotyping-by-sequencing (GBS) library, quality (molecular weight and purity) and quantity (concentration) of DNA are crucial and the

most significant technical issue. For instance, at the Institute of Genomic Diversity (IGD), Cornell University, there are stipulated requirements for GBS library preparation.

### *Agarose gel electrophoresis*

Yield and integrity of DNA can be determined using agarose gel electrophoresis (Fig. 7). Below are the standard procedures.

- Weigh 1% (1 g) agarose in 100ml 1XTBE buffer.
- Microwave to melt the agarose until a clear solution is observed.
- Cool the mix through continuous stirring and add 6 $\mu$ l of Ethidium bromide
- Transfer to an agarose gel apparatus with tray and appropriate comb inserted.
- Prepare 6 $\mu$ l mix consisting 3 $\mu$ l aliquots of DNA samples and 3 $\mu$ l of 2  $\times$  loading dye on a 96 well plate and briefly centrifuge (30 sec) at 1500 rpm.
- Load the mix and run the electrophoresis at constant voltage of 100 V for 1 hr.
- Also load lambda DNA standards of known concentration (100 ng/ $\mu$ l) on both ends of the gel for comparison of band intensity and degradation.
- Visualize and take gel picture under UV transilluminator by gel documentation with mounted camera (Fig. 7).

### *DNA quantity assessment*

In order to meet the DNA quantity requirement, 30  $\mu$ L of DNA at 30-100 ng/ $\mu$ L (as quantified by fluorometry), for GBS library construction at the Genomic Diversity Facility, we quantified all the samples using the multiple nanodrop. Samples with concentrations > 100 ng/ $\mu$ L will need to be diluted accordingly using the adjustment formula in Table 5, and then quantified again using the nanodrop. We usually discourage the submission of samples of low concentration because of poor quality and the inadequate reads that are often generated. Thus we re-extract any sample below 30 ng/ $\mu$ L, especially when the leaf tissue is still available and sufficient.

**Table 5: Comparing the nanodrop reading concentration and fluorometer concentration.**

Sample ID	Nanodrop reading	Equivalent fluorometer reading
TME 117	100 ng/ $\mu$ l	10 ng/ $\mu$ l
TMS961089A	1000 ng/ $\mu$ l	100 ng/ $\mu$ l

### *Restriction digestion*

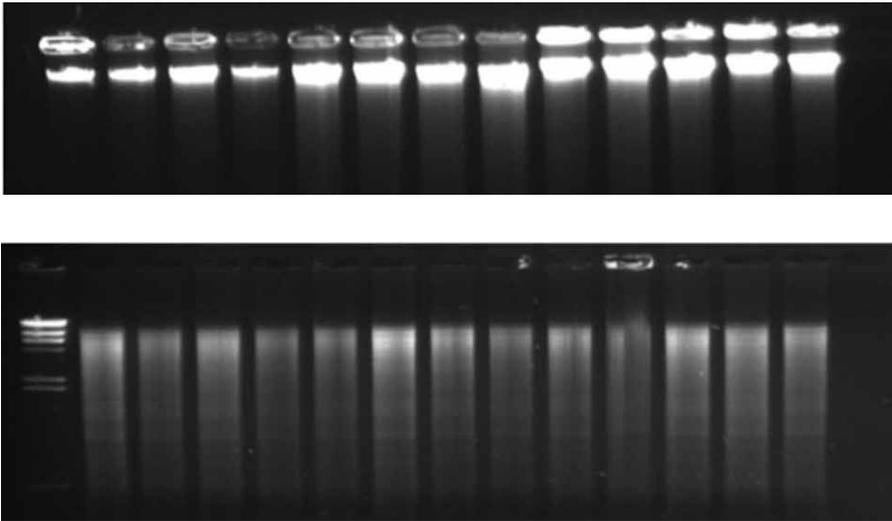
In order to further assess the quality of the extracted DNA, digestion is required with cheap restriction enzymes such as a non-methylation sensitive HindIII (a type II site-specific deoxyribonuclease restriction enzyme that cleaves the DNA palindromic

sequence AAGCTT). Hence below are the step-by-step procedures for restriction digestion of DNA samples with the above-mentioned restriction enzyme according to the manufacturer's protocol (Table 6).

- Dilute DNA samples to 200 ng/μL using a Tris-EDTA buffer (10mM Tris-HCl and 1mM EDTA at pH of 8.0).
- Print 3 μl of template DNA to a sterile 96-well plate and prepare a reaction mix for the digestion of the genomic DNA by adding 10 units of each restriction enzyme, 1 μl NE Buffer and make up to 10μl with sterile water.
- Gently mix all the tubes to avoid inactivation of the enzyme, centrifuge briefly, and incubate in the thermocycler at 37 °C for 3 hr.
- After incubation, mix 5 μl of digested DNA with 2 μl 5× loading buffer and run on a 1% agarose gel along with the uncut sample DNA (100 ng) and 500 ng of λ *HindIII* size standards at 100 volts for 1 hr.
- Visualize the gel images under UV trans-illuminator, annotate clearly using Microsoft PowerPoint and save in picture formats such as jpg, gif, tiff.
- Upload these gel images with limited size of 0.9MB on LIMS facility to allow the service provider to judge the quantity and quality of your sample DNA for GBS library construction.

**Table 6: Hind III digestion master mix.**

Reaction component	Volume (μl) ×1	Master mix (μl) ×100
Sterile deionized water	5	500
RE 10 × buffer	1	100
Restriction Enzyme ( <i>Hind III</i> ) (10 units/μl)	1	100
DNA (200 ng/μl)	3	
Total reaction volume	10	



**Figure 7: DNA quality and quantity assessment including test gel electrophoresis of DNA samples (top) and restriction enzyme digested for randomly selected DNA samples (bottom).**

### **Laboratory Information Management Systems (LIMS)**

A laboratory information management system (LIMS) is a software-based system with features that support the modern laboratory's operations. In relation to the CMS project we used the LIMS developed by the Genomic Diversity Facility (GDF) at Cornell University<sup>1</sup>. More information on the LIMS of the GDF is available on this link <https://slims.biotech.cornell.edu/default.aspx> under "User Guide".

### **DNA sample preparation and shipment for genotyping to other labs**

The following are important consideration during DNA sample preparation and shipment to other labs for outsourcing the genotyping work.

1. Transfer a total of 30  $\mu$ l DNA per sample to 96-well PCR plates. Make sure the blank well in each plate is blocked to ensure the well is devoid of liquid. This well is used as negative control for DNA sequencing.
2. Seal plates well with caps to avoid leaking and sample cross-contamination before centrifuging for about 1 min.
3. Just before packing, confirm that "blank" wells of corresponding plates are empty
4. Arrange plates in a stack (Fig. 8) and enclose inside the box containing cushioning materials to prevent plate damage during transit.
5. Ensure proper labeling (barcode label and distinct plate ID) for tracking and easy identification of different DNA sample plates.

NB: -Make sure a permanent marker is used for writing plate ID.

<sup>1</sup> Sequencing-based genotyping service can be offered by other service providers with each one of them using variants of LIMS system. Researchers should be able to find a suitable service provider for their projects depending on service, cost, turn-around time, and support.

6. Include all the necessary printed information such as the shipping address, plate ID number, and plate name and date as provided by the Genomic Diversity Facility, Cornell University.
7. Keep the samples in  $-20^{\circ}\text{C}$  until the final minutes before shipment.
8. Ship through couriers (e.g., DHL).
9. Plan the date of shipment in such a way that delivery at the other end will overlap on working days.
10. Shipment of samples at ambient temperature is highly recommended.



**Figure 8: Samples arranged in a stack (with barcode labels pasted on each plate) in addition to plate ID information.**

# Variety identification analysis

The pipeline for variety identification analysis has been developed and below is a summary of the details of the workflow from genotyping to variety identification including all the required and most important R-packages and other softwares.

## Genotyping by sequencing

Next-Generation Sequencing (NGS)-based genotyping procedures such as Genotyping-by-Sequencing (GBS) represent high-marker density approaches which can help to reveal the extent of genetic relatedness and genetic variation within and between cultivated species (Spindel et al. 2013). The GBS approach is based on reducing genome complexity with restriction enzymes coupled with multiplex NGS for the discovery of high-density single nucleotide polymorphism (SNP) markers (Elshire et al. 2011). The genome-wide molecular marker discovery, highly multiplexed genotyping, flexibility and low cost of GBS make it an excellent tool in plant genetics and breeding (Deschamps et al. 2012; Poland and Rife 2012).

Once quality DNA is available the next step in the GBS procedure is GBS library preparation by selecting and optimizing restriction digestion with restriction enzymes. For example, the ApeKI restriction enzyme (recognition site: G|CWCG) that produces less variable distributions of read depth and therefore a larger number of scorable SNPs in cassava was used in the study related to the CMS project. Depending on species and the number of markers desired, multiples of 96, 192, or 384-plex GBS libraries are required to be constructed for large-scale samples following the standard procedure (Elshire et al. 2011) and sequencing will follow using the Illumina HiSeq2500 or any other high throughput next generation sequencing platforms.

## SNP discovery and quality control of the SNP data

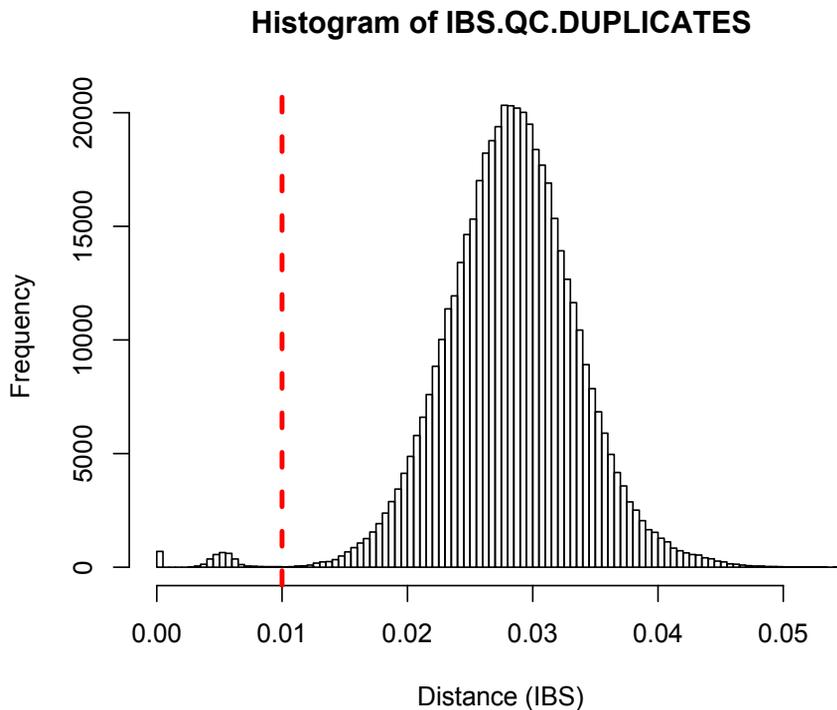
The raw read sequences of all the genotypes collected from different sources, including accessions in the reference library and duplicate samples should be processed through a TASSEL-GBS discovery pipeline developed using TASSEL 5.0 (Glaubitz et al. 2014). Perform a SNP calling based on TASSEL-GBS production pipeline by aligning the tags to the reference genome.

Quality control requires to be done by first calculating the missing (NA) data by SNP and individual followed by filtering with the missing genotype, missing individual, and minor allele frequency (MAF) using either TASSEL software or PLINK: Whole genome data analysis toolset version 1.9 (Purcell et al. 2007).

## Establishment of distance threshold to determine identical sets of genotypes

A few randomly selected samples (e.g., a total of 89 samples were used in the CMS study) are required to be genotyped in duplicates to determine a distance threshold between genotypes that will help in declaring a distance at which two genotypes or a set of genotypes are similar or distinct. A frequency distribution of distance (IBS) can be plotted and bimodal distribution of pairwise genetic distance obtained (see example in Figure 9).

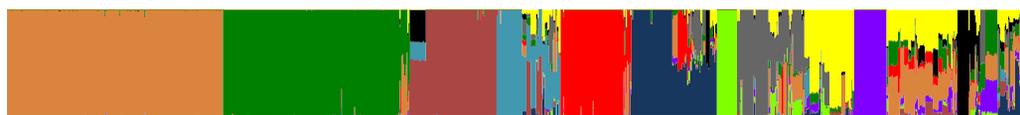
The bimodal distribution usually shows the frequency distribution of the data where one of the curves shows an artifact that could occur due to genotyping error. Filtering with MAF, missing individuals and missing genotypes will give the total genotyping rate and total variants (more robust SNP markers). The point between the bimodal distributions is therefore declared as a distance threshold where any pair of genotypes or set of genotypes below the point are identical. This “calibration principle” approach is taken because of the possibility of SNP genotyping errors resulting from miscalling some heterozygous SNPs with low sequencing read depth as homozygotes.



**Figure 9. Pairwise genetic distance (IBS) calculated using SNP markers. Note the red vertical line denoting the distance threshold below which two samples can be considered identical.**

## Cluster analysis

Once the quality control is performed and the distance threshold is determined, identify the varieties by computing distance-based hierarchical clustering, a pairwise genetic distance (identity-by-state, IBS) matrix calculated based on the robust SNP markers obtained after filtering in PLINK: Whole genome data analysis toolset version 1.9 (Purcell et al. 2007). A Ward's minimum variance hierarchical cluster dendrogram can then be built from the IBS matrix using the Analyses of Phylogenetics and Evolution (ape) package (Paradis et al. 2004) in R software (R Core Team 2013). The critical distance threshold determined can then be applied for the whole data and individuals belonging to the same cluster group below the threshold can be identified as a single genotype (See Figure 10).



**Figure 10. Example of cluster analysis showing genetically identical varieties. Figure adapted from Rabbi et al. 2015. (a) Hierarchical clustering (Ward's minimum variance method) dendrogram. The red dashed line represents the empirically determined distance threshold developed from comparison of duplicated library samples. A distance of 0.05 below which two individuals can be declared identical. (b, bottom) Individual ancestry estimated from ADMIXTURE analysis. Individuals are represented as thin vertical lines partitioned into segments corresponding to the inferred membership in  $K = 11$  genetic clusters as indicated by the colors. The roman numerals show groups of clonal individuals with predominant ancestry membership in each of the 11 clusters.**

## ADMIXTURE analysis

In line with the plant breeder's right, ancestry inferences can be useful in estimating the impacts resulting from the use of its improved germplasm by other programs (Morris and Heisey 2003). This is because improved germplasm often moves easily throughout the network of plant breeding systems, resulting in research spillover benefits.

ADMIXTURE software tool (Alexander et al. 2009) can be used for maximum likelihood estimation of individual ancestries from multi-locus SNP genotype datasets.

## Development of a reference library for variety identification

Developing a well-curated comprehensive reference library in collaboration with the breeding programs is very important for tracking genotypes of interest otherwise DNA fingerprinting alone can be used only to establish baseline data. The quality of the reference library (genotype traceability and comprehensiveness) determines "level of success" in variety identification. Moreover, a similar genotyping procedure should be applied for the development of the reference library and requires SNP calling and variety identification analysis by combining the raw sequences from the reference library with genotypes pending identification. Likewise, filtering for quality control by missing data and MAF can be done for the combined data. A comprehensive reference library was developed for cassava by IITA (Rabbi et al. 2015).

### **Matching samples to those in the reference library**

Individual accessions in the same cluster below the distance threshold can be considered as the same genotypes. If any of the genotypes from the reference library fall in the cluster of different individuals representing the same variety then they will be identified based on the genotypes from the reference library.

# References

- Alexander, D.H., J. Novembre, and K. Lange. 2009. Fast model-based estimation of ancestry in unrelated individuals. *Genome Research* 19(9):1655–1664.
- Dellaporta, S.L., J. Wood, and J.B. Hicks. 1983. A plant DNA miniprep protocol. Version II. *Plant Molecular Biology Reporter* 1: 19–21.
- Deschamps, S., V. Llaca, and G.D. May. 2012. Genotyping-by-Sequencing in Plants. *Biology* 1(3): 460–483.
- Elshire, R.J., J.C. Glaubitz, Q. Sun, J.A. Poland, and K. Kawamoto. 2011. A robust, simple Genotype-by-sequence (GBS) approach for high diversity species. *PLoS One* 6(5): e 19379. Doi:10.1371/journal.pone.0019379
- Morris, M.L. and P.W. Heisey. 2003. Estimating the benefits of plant breeding research: methodological issues and practical challenges. *Agricultural Economics* 29(3):241–252.
- Paradis, E., J. Claude, and K. Strimmer 2004. APE: Analyses of phylogenetics and evolution in R language. *Bioinformatics* 20: 289–290.
- Poland, J.A. and T.W. Rife. 2012. Genotyping-by-Sequencing for plant breeding and genetics. *Plant Genome* 5: 92–102.
- Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M.A. Ferreira, D. Bender, et al. 2007. PLINK: a tool set for whole-genome association and population-based linkage analyses. *American Journal of Human Genetics* 81: 559–575.
- Rabbi, I., M. Hamblin, M. Gedil, P. Kulakow, M. Ferguson, A.S. Ikpan, D. Ly, and J.L. Jannink. 2014. Genetic mapping using Genotyping-by-Sequencing in the clonally propagated cassava. *Crop Science* 54(4):1384–1396.
- Rabbi, I.Y., P.A. Kulakow, J.A. Manu-Aduening, A.A. Dankyi, I.Y. Asibuo, E.Y. Parkes, T. Abdoulaye, G. Girma, M.A. Gedil, P. Ramu, B. Reyes, and M.K. Maredia. 2015. Tracking crop varieties using genotyping-by-sequencing markers: a case study using cassava (*Manihot esculenta* Crantz). *BMC Genetics* 16: 115.
- R Core Team. 2013. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.
- Spindel, J., M. Wright, C. Chen, J. Cobb, J. Gage, S. Harrington, et al. 2013. Bridging the genotyping gap: using genotyping by sequencing (GBS) to add high-density SNP markers and new value to traditional bi-parental mapping and breeding populations. *Theoretical and Applied Genetics* 126: 2699–2716.

# Web resources

Barcode reader/scanner

<https://www.socketmobile.com/products/series-7/colorfuloverview/?languageRedirect=true>

Cornell University, Genomic Diversity Facility

<http://www.biotech.cornell.edu/brc/genomic-diversity-facility>

Garmin eTrex GPS device Manual

[http://static.garmincdn.com/pumac/eTrex\\_10-20-30\\_OM\\_EN.pdf](http://static.garmincdn.com/pumac/eTrex_10-20-30_OM_EN.pdf)

Garmin eTrex GPS device Tutorial videos

<http://www8.garmin.com/learningcenter/on-the-trail/etrex/>

DNA Plate application software

<http://wheatgenetics.org/research/12-research/software/21-dna-plate-app>

IDAutomation barcode software for barcode font

<http://www.idautomation.com/barcode-fonts/code-39/download.html>

IITA bioscience center

<http://bioscience.iita.org>

Permanent laser labels for containers

<http://www.phenixresearch.com>

Science laboratory safety (MSDS)

<https://www.sciencelab.com/msdsList.php>

Troubleshooting common problems in GPS devices

<http://www.wikihow.com/Troubleshoot-Common-Problems-with-a-Gps-Navigation-Unit>



