



# **Data Collection Guidelines**

Document Number: 30173

Document Revision: F

## Table of Contents

|   |    |
|---|----|
| Revision History .....                                | 4  |
| Data Collection Guidelines .....                      | 5  |
| Introduction .....                                    | 5  |
| Throughput Targets and Coverage .....                 | 5  |
| Coverage Recommendations for Different Analyses ..... | 7  |
| Resource Considerations .....                         | 10 |
| Example Scenarios .....                               | 10 |
| Technical Assistance .....                            | 17 |

## Legal Notice

---

### **For Research Use Only. Not for use in diagnostic procedures.**

This material is protected by United States Copyright Law and International Treaties. Unauthorized use of this material is prohibited. No part of the publication may be copied, reproduced, distributed, translated, reverse-engineered or transmitted in any form or by any media, or by any means, whether now known or unknown, without the express prior permission in writing from Bionano Genomics. Copying, under the law, includes translating into another language or format. The technical data contained herein is intended for ultimate destinations permitted by U.S. law. Diversion contrary to U. S. law prohibited. This publication represents the latest information available at the time of release. Due to continuous efforts to improve the product, technical changes may occur that are not reflected in this document. Bionano Genomics reserves the right to make changes in specifications and other information contained in this publication at any time and without prior notice. Please contact Bionano Genomics Customer Support for the latest information.

BIONANO GENOMICS DISCLAIMS ALL WARRANTIES WITH RESPECT TO THIS DOCUMENT, EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO THOSE OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. TO THE FULLEST EXTENT ALLOWED BY LAW, IN NO EVENT SHALL BIONANO GENOMICS BE LIABLE, WHETHER IN CONTRACT, TORT, WARRANTY, OR UNDER ANY STATUTE OR ON ANY OTHER BASIS FOR SPECIAL, INCIDENTAL, INDIRECT, PUNITIVE, MULTIPLE OR CONSEQUENTIAL DAMAGES IN CONNECTION WITH OR ARISING FROM THIS DOCUMENT, INCLUDING BUT NOT LIMITED TO THE USE THEREOF, WHETHER OR NOT FORESEEABLE AND WHETHER OR NOT BIONANO GENOMICS IS ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

### **Patents**

Products of Bionano Genomics® may be covered by one or more U.S. or foreign patents.

### **Trademarks**

The Bionano Genomics logo and names of Bionano Genomics products or services are registered trademarks or trademarks owned by Bionano Genomics in the United States and certain other countries.

Bionano Genomics®, Irys®, IrysView®, IrysChip®, IrysPrep®, IrysSolve®, Saphyr®, Saphyr Chip®, Bionano Access®, and Bionano EnFocus™ are trademarks of Bionano Genomics, Inc. All other trademarks are the sole property of their respective owners.

No license to use any trademarks of Bionano Genomics is given or implied. Users are not permitted to use these trademarks without the prior written consent of Bionano Genomics. The use of these trademarks or any other materials, except as permitted herein, is expressly prohibited and may be in violation of federal or other applicable laws.

© Copyright 2021 Bionano Genomics, Inc. All rights reserved.

## Revision History

---

| Revision | Release Date | Notes  |
|----------|--------------|--|
| A        | 5/9/2017     |  |
| B        | 4/5/2018     |  |
| C        | 5/16/2019    |  |
| D        | 10/16/2020   | Add guidance for RVA, FSHD. Add example scenarios  |
| E        | 11/04/2021   | Document updates for Solve 3.7 release <ul style="list-style-type: none"><li>• Clarify minSites <math>\geq 9</math> role in effective coverage calculations</li><li>• Add guidance for Fragile X</li><li>• Add file size table</li><li>• Update informatics report formats</li></ul> |
| F        | 01/20/2023   | Changed heading in column C of Table 4 and added new column. Also added text regarding sample quality.   |

## Data Collection Guidelines

---

### Introduction

The Saphyr® System collects data by drawing ultra-high molecular weight DNA into nanochannel arrays, imaging the molecules in linear confinement, and detecting molecule and label positions from the images. Once a batch of molecules in the nanochannels have been imaged, the Saphyr clears the array and loads in new molecules. The system repeatedly cycles through this process (each cycle is one “scan”). The greater the number of scans that accumulate, the rawer data collects for a given dataset. Operational details are covered in the **Saphyr System User Guide - 30247**.

Per flowcell limits on maximum accumulated output and runtime vary by chip part number and Saphyr Instrument Control Software (ICS) version. Within those limits, it is up to the user to determine how much data to collect. This document will provide guidance and discuss data collection strategies for various downstream applications.

The key factors for success in using the Saphyr System include:

- Design your experiments to meet research goals.
- Isolate high molecular weight gDNA from sample sources, then efficiently label at enzyme-specific recognition sequence.
- Collect sufficient amount of high quality data for downstream data analysis.

The Saphyr System collects whole genome imaging data. As such, throughput and coverage targets discussed here are genome-wide in scope, regardless of any more focused downstream application. Furthermore, the guidance here assumes that data is of reasonably good quality. For a discussion of raw data quality, please refer to **Bionano Access Dashboard Guidelines – 30304** and/or **Molecule Quality Report Guidelines – 30223**. For guidance on troubleshooting sample prep or quality from a particular sample type, please refer to the appropriate Bionano Prep™ documentation.

### Throughput Targets and Coverage

Prior to collecting Saphyr data, the user enters chip and flowcell information in Bionano Access. A user is prompted to specify the Throughput Target (in Gbp) for each flowcell. This throughput target should reflect an appropriate fold coverage (X) of the haploid genome size investigated, specific to the intended application.

Conceptually, coverage refers to the average depth of molecules representing a given locus in the genome. It can be calculated from the predicted haploid genome size and the amount of input data.

$$\text{Coverage (X)} = \frac{\text{Total DNA}}{\text{Genome size}}$$

For example, 310 Gbp Total DNA Throughput of a sample with a 3.1 Gbp haploid genome size equals 100X Coverage (310 Gbp ÷ 3.1 Gbp = 100).

In practice, not every molecule in a raw dataset adequately represents a locus. Similarly, reference assemblies are not perfect approximations of the genome. It is useful to define *effective coverage of the reference*, which uses alignment (*i.e.* Map Rate) and genome size as in reference assembly as parameters.

$$\text{Effective coverage of reference (X)} = \frac{\text{Total DNA} \times [\text{Map Rate}]}{\text{Reference size}}$$

Throughput and Coverage are further refined into specific terms found throughout Bionano software, outlined in **Table 1**.

**Table 1. Glossary of Terms relating Throughput and Coverage**

| Term   | Meaning   | Source  | Comment  |
|--|---|---|--|
| <b>Pre-Analysis metrics: prior to secondary bioinformatic analysis</b>               |   |   |  |
| Raw coverage (X)   | [Total DNA (>=150 kbp)] ÷ [haploid genome size]                             |   | Useful baseline when reference and/or data quality unclear         |
| Throughput Target  | Minimum total flowcell throughput (≥150 kbp) that Saphyr targets to collect | Bionano Access; Saphyr Instrument Control Software (Saphyr ICS) |  |
| Total DNA (>= 20 kbp)  | Cumulative length of DNA molecules that are ≥ 20 kbp                        | Dashboard; Molecule Quality Report                              | Raw molecules detectable by Saphyr                                 |
| Total DNA (>=150 kbp)  | Cumulative length of DNA molecules that are ≥ 150 kbp                       | Dashboard; Molecule Quality Report                              | Counts toward Throughput Target                                    |
| Est Effective coverage >= 150 kbp*   | [Total DNA (>=150 kbp)] × [Map Rate %]* ÷ [reference size]                  | Dashboard   | Approximates effective coverage of reference (X) prior to analysis |
| Total DNA (>= 150 kbp & minSites >= 9)   | Cumulative length of DNA molecules that are ≥ 150 kbp and have ≥ 9 labels   | Molecule Quality Report   | Pass filters for downstream analyses                               |
| Effective coverage*  | [Total DNA (>= 150 kbp & minSites >= 9)] × [Map Rate %]* ÷ [reference size] | Molecule Quality Report   | Approximates effective coverage of reference (X) prior to analysis |
| <b>Post-Analysis metrics: output from secondary bioinformatic analysis pipelines</b> |   |   |  |

|                                      |  |   |  |
|--------------------------------------|--|---|--|
| Effective coverage of reference (X)* | $[\text{Total DNA aligned to the reference in pipeline}]^* \div [\text{reference size}]$ | De Novo Assembly Report; Rare Variant Analysis Report |  |
| Effective coverage of assembly (X)   | $[\text{Total DNA aligned to the assembly in pipeline}] \div [\text{assembly size}]$     | De Novo Assembly Report                               |  |

\*Calculation of *Map Rate* and *Effective Coverages* relies on quality and completeness of the provided reference, as well as sequence similarity between reference and sample. Human genome references are considered relatively complete, so these metrics are meaningful and can be clearly benchmarked. With non-human organisms or custom references, molecule-reference alignment can be hindered by a myriad of factors, which may underestimate coverage and data quality. Please see **Assembly Report Guidelines (#30255)** for additional details.

In human samples, the effective coverage provided pre-analysis is generally higher than the effective coverage of reference (X) in the pipeline results. This stems from different stringencies in alignment parameters between pre-analysis mapping and eventually running the pipeline. Effective coverage recommendations outlined here refer to pre-analysis metric shown in the Molecule Quality Report.

## Coverage Recommendations for Different Analyses

Coverage recommendations vary considerably by downstream application. As of Bionano Solve 3.7, there are pipelines for Structural Variant and Copy Number Analysis (De Novo Assembly and Rare Variant Analysis), Facioscapulohumeral muscular dystrophy (FSHD) Analysis, Fragile X Analysis, and Hybrid Scaffolding. **Tables 2-5** provide a summary of minimum coverage recommendation for Direct Label and Stain (DLS) data for each application. For a more detailed view of rationale (such as SV sensitivity or CN calling performance at different coverages) please refer to the Theory of Operations guide appropriate to the analysis at [www.bionanogenomics.com/support/](http://www.bionanogenomics.com/support/).

Beginning in Solve 3.4, Bionano offers two distinct pipelines for whole genome calling of structural variants, variant annotation, and copy number variants: the *de novo* assembly pipeline and rare variant analysis. The *de novo* assembly pipeline is intended for downstream analysis of constitutional, germline, and “typically” diploid SVs. It requires a modest amount of coverage (**Table 2**) to assemble large genome maps *de novo*, then aligns maps to a provided reference. Genome structure is elucidated, and SVs called are generally assumed to be homozygous or heterozygous.

**Table 2. De Novo Assembly recommended minimum input, targeting human SV analysis**

| Pipeline                | Minimum Pre-Analysis Effective Coverage (X) | hg19 / hg38 length (Gbp) | Example Map Rate    | Min. DNA (Gbp)      | Post-Analysis Effective Coverage of Reference(X) |
|-------------------------|---|--------------------------|---------------------|---------------------|--|
| <b>De Novo Assembly</b> | <b>80</b>                                   | <b>× 3.1</b>             | <b>÷ <u>78%</u></b> | <b>= <u>318</u></b> | <b>≈70</b>                                       |

Bionano Genomics has observed very high sensitivity and confidence (PPV) at 80X effective coverage pre-analysis for heterozygous and homozygous SVs in non-mosaic cases. 80X is a recommended minimum, which translates to 320Gbp of DNA collected by Saphyr for a typical run with good molecule metrics. More coverage will yield small additional sensitivity for heterozygous and homozygous SVs with diminishing returns beyond 120X for non-mosaic cases. If a different effective coverage is desired, simply substitute a new coverage target in the calculation table above.

In contrast to *de novo* assembly, Rare Variant Analysis is intended for analysis of cancer, somatic, or otherwise heterogenous samples. It leverages deeper input coverage (**Table 3**) to target SVs at low allelic fractions. The pipeline directly aligns molecules to hg19 or hg38 to capture SVs with a relatively low fraction of molecule support.

**Table 3. Rare Variant Analysis recommended minimum input, for targeting ≥5% Variant Allele Fraction**

| Pipeline              | Minimum Pre-Analysis Effective Coverage (X) | hg19 / hg38 length (Gbp) | Example Map Rate | Min. DNA (Gbp) | Post-Analysis Effective Coverage of Reference(X) |
|-----------------------|---|--------------------------|------------------|----------------|--|
| Rare Variant Analysis | 340 <sup>†</sup>                            | × 3.1                    | ÷ <u>78%</u>     | = <u>1351</u>  | ≈300 <sup>†</sup>                                |

<sup>†</sup>Bionano Genomics has observed at least 90% sensitivity to structural variants at 5% variant allele fraction (VAF), with datasets at 300X effective coverage of reference post-analysis. Data input of lower coverage is acceptable but VAF sensitivity will decrease. This performance is shown in detail in the **Theory Of Operation – Structural Variant Calling - 30110**.

To ensure adequate effective coverage, we generally recommend a Throughput Target of 1500 Gbp for datasets intended for Rare Variant Analysis targeting ≥5% Variant Allele Fraction. The pipeline has been validated with datasets of up to 5000 Gbp Total DNA (≥150 kbp). This very high degree of coverage can improve sensitivity beyond typical low variant allele frequency applications (false discovery rate at 5000 Gbp has not been extensively examined). If a different effective coverage is desired, simply substitute a new coverage target in the calculation table above.

The Bionano EnFocus™ FSHD Analysis and EnFocus™ Fragile X Analysis use targeted assembly approaches. The pipelines define minimum acceptable criteria for data input and performs a PASS/FAIL assessment of molecule quality. The whole genome coverage recommendation provided is a safe minimum shown to confidently resolve targeted regions. Criteria are shown in **Table 4**. Of note, the post-auto noise correction effective coverage will likely be affected by the quality of the sample.

**Table 4. EnFocus™ FSHD Analysis and EnFocus™ Fragile X Analysis Minimum Data Criteria**



| Analysis                  | Pipeline                           | Post Auto-noise Correction Effective Coverage (X) | Map Rate    | N50 (>= 150 kbp) | Recommended Min. DNA (Gbp) |
|---------------------------|------------------------------------|---|-------------|------------------|----------------------------|
| <b>FSHD Analysis</b>      | <b>EnFocus™ FSHD Analysis</b>      | <b>≥75</b>  | <b>≥70%</b> | <b>≥200 kbp</b>  | <b>=400</b>                |
| <b>Fragile X Analysis</b> | <b>EnFocus™ Fragile X Analysis</b> | <b>≥75</b>  | <b>≥70%</b> | <b>≥200 kbp</b>  | <b>=400</b>                |

Finally, the hybrid scaffolding pipeline utilizes one or two *de novo* assembly output(s) to scaffold NGS contigs/scaffolds. It is normally used with plant and animal haploid NGS assemblies, frequently of genomes that are otherwise poorly characterized. The first step in a hybrid scaffolding workflow is collecting data to support *de novo* assembly to be used as input alongside NGS data. The assembly coverage recommendation is what is required to build contiguous and accurate consensus genome maps.

**Table 5. De Novo Assembly recommended minimum input targeting downstream use with Hybrid Scaffolding, using Saphyr data and reference derived from sample matched haploid NGS with reasonable completeness and contiguity**

| Pipeline                | Minimum Pre-Analysis Effective Coverage (X)* | Haploid genome size (Gbp) | Example Map Rate*   | Min. DNA (Gbp)      | Post-Analysis Effective Coverage of Reference (X) |
|-------------------------|--|---------------------------|---------------------|---------------------|---|
| <b>De Novo Assembly</b> | <b>80</b>                                    | <b>× <u>1.1</u></b>       | <b>÷ <u>55%</u></b> | <b>= <u>160</u></b> | <b>≥70</b>  |

\*Calculation of *Map Rate* and *Effective Coverage of the Reference (X)* relies on quality and completeness of the provided reference, as well as sequence similarity between reference and sample. For example, the *Map Rate* shown above (55%) would be considered low for human data against the highly complete hg19. It may be sufficient in a non-model system. Following calculation in **Table 5**, reduced map rates could cause an overestimation of Min DNA (Gbp) needed to meet target pre-analysis effective coverage. See example scenario 2 below.

If reference quality is poor or none is available, 100-120X raw haploid coverage can serve as minimum input target.

Genome researchers have enjoyed significant size and accuracy gains from scaffolding NGS contigs/scaffolds with *de novo* assembled Bionano maps. 80X is a recommended minimum for this application. When using a Direct Label and Stain (DLS) approach, effective coverage up to and beyond 100X has shown improved map contiguities for some plants and animals. Expect pipeline output to be highly genome specific. If a different effective coverage is desired, simply substitute a new coverage target in the calculation table above.

## Resource Considerations

Collecting and/or analyzing high volumes of data incurs costs. These costs include *i.* time a chip spends occupying the Saphyr System, *ii.* data storage consumption, and *iii.* token (for Compute On Demand operations) and runtime costs of completing analysis pipelines. Runtime estimates provided by Bionano are based on our recommended coverage levels.

- As measured by running the Bionano control sample, throughput of G2.3 chips in a Saphyr System can support collecting three datasets of  $\geq 320$  Gbp (fit for human *de novo* assembly; **Table 2**) every 8 hours. For targeting  $\geq 1500$  Gbp (Rare Variant Analysis; **Table 3**), data collection for three samples generally can take around 24-36 hours. Applications requiring higher data collection will take longer. For example, Bionano has observed collection of full 5000 Gbp datasets take up to 96 hours. Any time dedicated to collecting a data volume beyond necessary coverage, poses some availability cost against collecting the next datasets.
- Larger volume datasets require larger files which consume storage space faster. Large file sizes may create difficulties with file transfer and data backup. Please see **Table 6** for an approximation of file sizes.
- Submitting higher volume datasets for bioinformatic analyses increases the computational burden. Higher data inputs increase job runtime, imposing an availability cost on analysis goals. For analyses through Compute On Demand, higher coverage than necessary needlessly increases token cost.

If a dataset is larger than necessary for a desired application, we recommend to downsample it in Bionano Access (see **Bionano Access® Software User Guide - 30142**, and example scenario 1b).

**Table 6. Representative sizes for different file types and volumes.**

| Output Type                     | Extension | Total DNA ( $\geq 150$ kbp) input | Size* (compressed) |
|---------------------------------|-----------|-----------------------------------|--------------------|
| Molecules*                      | .bnx.gz   | 400 Gbp                           | 1.1 GB             |
|                                 |           | 800 Gbp                           | 1.6 GB             |
|                                 |           | 1500 Gbp                          | 3.6 GB             |
| Annotated De Novo Assembly      | .zip      | 400 Gbp                           | 2.4 GB             |
|                                 |           | 800 Gbp                           | 4.5 GB             |
| Annotated Rare Variant Analysis | .zip      | 1500 Gbp                          | 3.8 GB             |
| EnFocus™ FSHD                   | .zip      | 400 Gbp                           | 140 MB             |
| EnFocus™ Fragile X              | .zip      | 400 Gbp                           | 133 MB             |

\*Sizes assume good raw data quality, with reasonably low proportion of molecules  $< 150$  kbp. Datasets with excessive DNA fragmentation may inflate file sizes.

## Example Scenarios

### Scenario 1a.

**Research Goal:** Rare Variant Analysis, targeting  $\geq 5\%$  variant allele fraction

In this scenario, a cancer researcher wishes to interrogate a sample for structural variants at  $\geq 5\%$  variant allele fraction. The researcher enters deidentified sample information into Bionano Access and sets the Throughput Target (Gbp) to 1500. The flowcell completes in 22 hours, and generates the following Molecule Quality Report (MQR):

## MQR Report Details

| label   | value                      | description   |
|---|----------------------------|---|
| Reference   | hg38_DLE1_0kb_0labels.cmap | Name of the reference genome this sample was aligned to.                                    |
| Reference Length                                    | 3,088,269,832 bp           | Total length of reference sequence  |
| Enzyme  | DLE-1                      | Name of the enzyme used in this sample.   |
| Site  | CTTAAG                     | Recognition sequence of the enzyme used.  |
| Maximum molecule length                             | 2.2 Mbp                    | The longest molecule detected during the chip run.  |
| N50 ( $\geq 20$ kbp)                                | 232.21 kbp                 | N50 of the molecules that are 20kbp or longer)  |
| Total DNA ( $\geq 20$ kbp)                          | 2,213.93 Gbp               | Total amount of DNA from molecules that are 20 kbp or longer                                |
| N50 ( $\geq 150$ kbp)                               | 299.18 kbp                 | N50 of DNA molecules that are 150kbp or longer  |
| Total DNA ( $\geq 150$ kbp)                         | 1,527.11 Gbp               | Total amount of DNA from molecules that are 150kbp or longer                                |
| N50 ( $\geq 150$ kbp and min sites $\geq 9$ )       | 301.06 kbp                 | Same as other N50 fields, but molecules must have at least 9 labels                         |
| Total DNA ( $\geq 150$ kbp and min sites $\geq 9$ ) | 1,483.58 Gbp               | Same as other Total DNA fields, but molecules must have at least 9 labels                   |
| Map rate  | 82.3 %                     | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage                                  | 395.40                     | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |
| Average label density ( $\geq 150$ kbp)             | 15.93 /100kbp              | Average number of labels per 100 kbp for the molecules that are 150kbp or longer            |
| Site SD   | 0.11                       | Constant term in sizing error relative to reference   |
| Relative SD   | 0.021                      | Quadratic term in sizing error relative to reference  |
| Scaling SD  | 0                          | Linear term in sizing error relative to reference   |
| integrity_num                                       | 0.11                       |   |

| label                         | value      | description   |
|-------------------------------|------------|---|
| Negative label variance (NLV) | 10.7       | Percentage of reference labels absent in molecules  |
| Base pairs per pixel          | 487.7      | Calculated base pairs per pixel in the alignment by comparing molecules to the reference. |
| Label color                   | BNGFLGR001 | Label color used for detection.   |
| version                       | 1          |   |
| Positive label variance (PLV) | 6.78       | Percentage of labels absent in reference  |

Referring to **Table 3**, she notes that pre-analysis effective coverage has reached the minimum for her application. The research team proceeds with Rare Variant Analysis.

**Scenario 1b.**

**Research Goal:** *De Novo* Assembly with dataset from 1a

Another researcher from the lab accesses the data and would like to assemble germline structure of the genome. The researcher decides to separately analyze the raw data with the *De Novo* Assembly pipeline. Based on the original MQR data from Scenario 1a, the pre-analysis effective coverage is much higher than necessary for the application.

To save on cost and runtime, the researcher selects to assemble a random downsample targeting 80X. Bionano Access enables her to select an 80X downsample during her job launch. Those same calculations from **Table 2** are:

$$[ 80X \text{ pre-analysis coverage } ] \times [ 3.1 \text{ hg38 length} ] \div [ 0.823 \text{ Map Rate } ] = [ \mathbf{301} \text{ Total DNA (} \geq 150 \text{ kbp \& minSites } \geq 9) ]$$

The researcher proceeds with *De Novo* Assembly.

**Scenario 2.**

**Research Goal:** *De Novo* Assembly to support Hybrid Scaffolding a small plant genome

An agricultural center wants to improve the reference genome of a novel plant species. The NGS assembly that will be scaffolded is not available at the time Bionano data is collected. However, reference sequences are available for related plants. The most similar one is chosen for *In Silico* Digestion (see **Bionano Access® Software User Guide – 30142**). The Files Summary produced is assessed for site density with DLE-1:

Channel 1 site density (sites/100kbp): 26.0  
Channel 1 estimated label density (labels/100kbp) for Saphyr: 19.2

Channel 1 estimated label density (labels/100kbp) for Irys: 16.8

As site density falls within recommended 8 – 31 / 100kbp range for direct labeling - the cmap is added to references in Access, and a technician labels plant DNA with DLE-1.

The researchers expect a haploid genome size of 850 Mbp and decide to target 110X raw coverage. They enter sample information into Bionano Access and set the Throughput Target (Gbp) to 94. The flowcell completes in five hours, and generates the following Molecule Quality Report (MQR):

## MQR Report Details

| label                                    | value                               | description   |
|--|-------------------------------------|---|
| Reference                                | similar-plant_DLE1_0kb_0labels.cmap | Name of the reference genome this sample was aligned to.                                    |
| Reference Length                         | 850,399,333 bp                      | Total length of reference sequence  |
| Enzyme                                   | DLE-1                               | Name of the enzyme used in this sample.   |
| Site                                     | CTTAAG                              | Recognition sequence of the enzyme used.  |
| Maximum molecule length                  | 1.98 Mbp                            | The longest molecule detected during the chip run.  |
| N50 (>= 20 kbp)                          | 124.44 kbp                          | N50 of the molecules that are 20kbp or longer)  |
| Total DNA (>= 20kbp)                     | 225.4 Gbp                           | Total amount of DNA from molecules that are 20 kbp or longer                                |
| N50 (>= 150kbp)                          | 207.42 kbp                          | N50 of DNA molecules that are 150kbp or longer  |
| Total DNA (>= 150kbp)                    | 94.52 Gbp                           | Total amount of DNA from molecules that are 150kbp or longer                                |
| N50 (>= 150kbp and min sites >=9)        | 211.75 kbp                          | Same as other N50 fields, but molecules must have at least 9 labels                         |
| Total DNA (>= 150kbp and min sites >= 9) | 94.0 Gbp                            | Same as other Total DNA fields, but molecules must have at least 9 labels                   |
| Map rate                                 | 29.9 %                              | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage                       | 32.85                               | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |
| Average label density (>= 150kbp)        | 18.19 /100kbp                       | Average number of labels per 100 kbp for the molecules that are 150kbp or longer            |
| Site SD                                  | 0.1                                 | Constant term in sizing error relative to reference   |
| Relative SD                              | 0.024                               | Quadratic term in sizing error relative to reference  |
| Scaling SD                               | 0                                   | Linear term in sizing error relative to reference   |

| label                         | value      | description  |
|-------------------------------|------------|--|
| integrity_num                 | 0.1        |  |
| Negative label variance (NLV) | 23.4       | Percentage of reference labels absent in molecules                     |
| Base pairs per pixel          | 500.4      | Calculated base pairs per pixel in the alignment by comparing molecule |
| Label color                   | BNGFLGR001 | Label color used for detection.  |
| version                       | 1          |  |
| Positive label variance (PLV) | 18.33      | Percentage of labels absent in reference                               |

The map rate of the sample to this reference is poor (<30%), with positive and negative label variance both rather high. The label density is within ~15% of expectation. These factors suggest labeled DNA quality may be relatively good, with sample-reference mismatch causing low map rate and coverage estimate problems. Sequence dissimilarity between sample and reference, and/or poor reference assembly contiguity, can cause these issues in non-model organisms.

The operator may continue to collect data in the flowcell or begin analysis with what he has collected. To account for the uncertainty, the operator decides to do both. A De Novo Assembly is launched on the dataset; here is output from the assembly report:

## Molecules aligned to the reference

| label                             | value  | description   |
|-----------------------------------|--------|---|
| Total number of molecules aligned | 310769 | The number of molecules after filtering ( $\geq 150$ kbp) that align either to the reference file (.cmap), (e.g. GRCh37 or GRCh38) or the assembly. |
| Fraction of molecules aligned     | 0.3    | The proportion of filtered molecules that align to the reference or consensus genome maps (assembly only).  |
| Effective coverage of reference   | 32.92  | The total length of molecules divided by the length of the reference or consensus assembled maps after de novo assembly.                            |
| Average confidence                | 21.1   | The average alignment score for all the molecules that align to the reference or assembly.  |

...

## Molecules aligned to the assembly

| label                             | value  | description  |
|-----------------------------------|--------|--|
| Total number of molecules aligned | 857354 | The number of molecules after filtering ( $\geq 150$ kbp) that align either to the reference file (.cmap), (e.g., GRCh37 or GRCh38) or the assembly. |
| Effective coverage of assembly    | 82.489 | The total length of molecules divided by the length of the reference or consensus assembled maps after de novo assembly.                             |
| Average confidence                | 36.2   | The average alignment score for all the molecules that align to the reference or assembly.   |

The researchers refer to **Assembly Report Guidelines – 30255** to evaluate their result and observe that molecule alignment against the assembly is much improved over alignment to the reference. Also, the assembly size was consistent with their expectation. They note that this particular reference is a poor guide for the sample’s map rate and pre-analysis effective coverage. The completed assembly looks adequate, and more data is incoming from the running flowcell if added coverage is desired. In the absence of a good reference, they may also consider running the de novo assembly with the pre-assembly option enabled.

### Scenario 3.

**Research Goal:** *De Novo* Assembly in human datasets to support *Variant Annotation: trio*. One dataset map rate is lower than expected.

In this example, a researcher wants to analyze germline SVs in a child with apparent genetic disease. The researcher plans a trio analysis of the proband and both parents. The lab draws three blood samples, which are labeled with DLE-1. The researcher adds a G2.3 chip to the experiment and populates the three flowcells with deidentified mother/father/proband information. Throughput Target (Gbp) is set to 320 for each. The chip completes in six hours, and generates the following three map rates and pre-analysis coverages:

...

-----

| label              | value | description   |
|--------------------|-------|---|
| Map rate           | 85.8% | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage | 92.52 | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |

...

| label              | value | description   |
|--------------------|-------|---|
| Map rate           | 83.3% | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage | 89.9  | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |

...

---

| label              | value | description   |
|--------------------|-------|---|
| Map rate           | 71.9% | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage | 75.21 | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |

...

The Saphyr operator observes that one samples finished below the recommended 80X effective coverage minimum. Its map rate is somewhat lower than expectation, which reduces input coverage as per calculation in **Table 2**. During data collection, its map rate stabilized lower than the other two samples.

The Throughput Target was adjusted accordingly. The flowcell was rehydrated and run by itself to a new total Throughput Target (Gbp) of 400. The additional runtime completes in 90 minutes, and the run generates the following total dataset map rate and pre-analysis coverage:

...

---

| label              | value | description   |
|--------------------|-------|---|
| Map rate           | 71.8% | Percentage of molecules that are 150kbp or longer mapped to the reference                   |
| Effective coverage | 93.7  | Total amount of aligned DNA divided by the size of the reference genome times the map rate. |

...

With all three datasets now  $\geq 80X$  effective pre-analysis coverage, the research team is satisfied all SVs annotated in the trio will have been generated with comparable sensitivity. They begin analyses with the *De Novo* Assembly and Variant Annotation Pipelines.



## Technical Assistance

---

For technical assistance, contact Bionano Genomics Technical Support.

You can retrieve documentation on Bionano products, SDS's, certificates of analysis, frequently asked questions, and other related documents from the Support website or by request through e-mail and telephone.

| Type    | Contact   |
|---------|---|
| Email   | <b>support@bionanogenomics.com</b>  |
| Phone   | <b>Hours of Operation:</b><br><br><b>Monday through Friday, 9:00 a.m. to 5:00 p.m., PST</b><br><br><b>US: +1 (858) 888-7663</b> |
| Website | <b><a href="http://www.bionanogenomics.com/support">www.bionanogenomics.com/support</a></b>                                     |