

Generating Accurate and Contiguous *De Novo* Genome Assemblies Using Hybrid Scaffolding

Bionano Optical Mapping Reveals True Long-Range Structure of the Genome while Reducing Sequencing Costs

Generating high-quality finished genomes remains challenging. Accurate identification of structural variations with minimal gaps is difficult or impossible using sequencing technologies alone.

The genomes of a majority of higher organisms are highly repetitive, with two thirds of human and most mammalian genomes comprised of repeats, and many plant genomes have even higher repeat content. Sequencing usually fails to span repeat arrays or disambiguate different copies of interspersed repeats that are not spanned. These failures can limit contig length and introduce chimeric joins and other assembly errors.

The widespread use of next-generation sequencing (NGS) has led to an accumulation of incomplete assemblies that contain large numbers of contigs but limited long-range information. NGS technology is based on fragmenting DNA molecules, reading just hundred(s) of base pairs, and then using algorithms to reassemble these fragments.

The introduction of long-read sequencing has led to improved assembly contiguity and accuracy, but can be time-consuming and expensive, especially when deep coverage or spanning of long tandem repeats is required. Read lengths are still limited to tens of kilobase pairs with highly fragmented genome assemblies.

Synthetic long-read technologies like that of 10x Genomics have similar limitations as NGS. Using a barcoding method to link short reads retains some mid-range structural information improving the contiguity of some NGS assemblies. However, synthetic long reads are still plagued by the same challenges inherent to NGS technology. These challenges include failure to disambiguate interspersed repeat units and correctly assemble and size long repeat arrays; assembly gaps due to incomplete coverage and GC bias in PCR amplification and sequencing; lack of long-range structural information caused by fragmenting of DNA and reads being too short to span and correctly resolve larger structural variations.

Only extremely long, megabase scale molecules provide accurate structure of the genome. Bionano optical mapping images long DNA molecules in their native state, while preserving long range genomic structural information. Structural variations are observed instead of algorithmically inferred as in sequencing approaches. These long-labeled molecules are *de novo* assembled into physical maps spanning the whole genome. The resulting label patterns in the map can be used for anchoring NGS contigs and detecting structural variants.

Megabase size molecules of genomic DNA are labeled at a specific 6-7 bp site, linearized and uniformly stretched in high density NanoChannel arrays, and imaged on the Saphyr™ instrument. For labeling, there are two distinct approaches. Traditionally, Bionano mapping used a nicking endonuclease to nick the sequence motif, followed by Labeling, Repair, and Stain (NLRS). This process is highly robust and specific, but it introduces systematic double-stranded breaks that limit the contiguity of Bionano maps.

The Direct Label and Stain method (DLS) does not nick the DNA, eliminating systematic molecule breaks. The DLS protocol consists of a single enzymatic labeling reaction, followed by cleaning and staining. DLS protocol is more streamlined since there is no need to repair DNA.

The label patterns generated by NLRS or DLS allow each long molecule to be uniquely identified and aligned. Using pairwise alignment of the single molecules, consensus genome maps are constructed, refined, extended and merged. DLS genome maps are up to 50-fold longer than NLRS maps, improving visualization of genome structure and creating the most contiguous and accurate assemblies. Chromosome arms and full chromosomes are often assembled in single maps. Multiple sets of genome maps can be created using different labeling enzymes and combined to generate broader coverage and higher label density.

| Sample | Molecule N50 > 150 Kbp (Kbp) | Bionano Map N50 (Mbp) |
|-------------------|---------------------------------|--------------------------|
| NA12878 | 293 | 55.9 |
| Human Fresh Blood | 307 | 56.9 |
| Bionano Maize B73 | 260 | 100.0 |
| Durum Wheat | 364 | 13.0 |
| Farro | 300 | 32.7 |
| Strawberry | 241 | 13.3 |
| Kakapo | 247 | 69.3 |
| Hummingbird | 310 | 38.7 |
| Blackbird | 243 | 21.6 |
| Fish | 245 | 22.3 |
| Ferret | 262 | 66.1 |
| Pig | 335 | 65.2 |
| Soybean | 246 | 23.0 |
| Brassica | 270 | 12.4 |
| Mouse | 280 | 101.0 |

Table 1:

Organisms de novo assembled using Bionano direct label stain chemistry (DLS). De novo assemblies often cover whole chromosome arms, only broken at centromeres and other low complexity regions which are longer than molecule length.

Hybrid Scaffold Construction

The *de novo* Bionano optical maps can be integrated with a sequence assembly to order and orient sequence fragments, identify and correct potential chimeric joins in the sequence assembly, and estimate the gap size between adjacent sequences. In order to do so, the Bionano Solve™ software imports the assembly and identifies label sites in the sequence based on the labeling specific recognition sequence motifs.

These *in silico* maps for the sequence contigs are then aligned to the *de novo* Bionano optical maps. Conflicts between the two are identified and resolved, and hybrid scaffolds are generated in which sequence maps are used to bridge Bionano maps and vice versa. Finally, the sequence assembly corresponding to this hybrid scaffold is generated and exported as FASTA and AGP files.

The pipeline is fully integrated with Bionano Access™ which provides a convenient interface for running the Hybrid Scaffold pipeline and viewing scaffolding results (Figure 1).

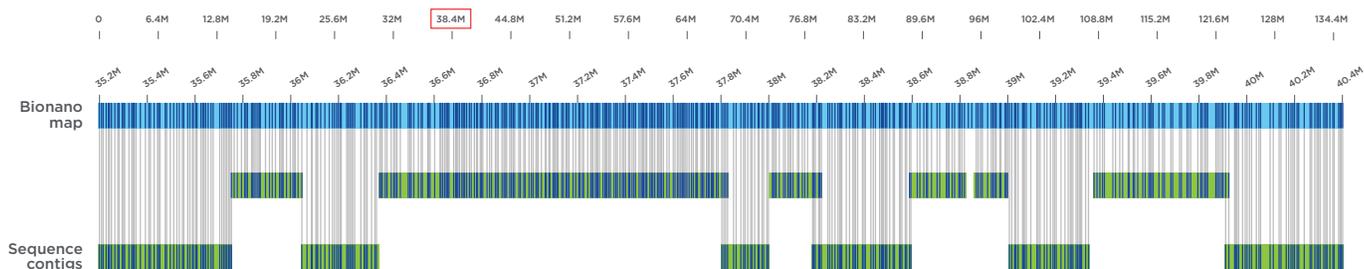


Figure 1: Combining Bionano maps with sequence assemblies. Sequence contigs are anchored and oriented using the de novo generated Bionano maps.

Contiguity and Completeness

The hybrid scaffolding process reduces thousands of contigs found in the initial NGS assembly to a handful of scaffolds, improving assembly accuracy and quality while reducing the need for deep sequencing coverage.

The hybrid scaffolding approach can yield significant improvements in contiguity across various species, as expressed by the assembly N50 values, as seen in Table 2. We created hybrid scaffolds for three genomes (maize, kakapo, and blackbird) which improved

contiguity by as much as 84-fold. The Bionano Solve pipeline makes near-complete use of the available input assemblies, taking into account 95-99.5% of the total length of the sequence (Table 2).

Bionano hybrid scaffolding is agnostic to the sequence technology used. Recent publications featured scaffolded assemblies using Illumina sequencing only¹, PacBio only², 10x Genomics assemblies³, nanopore sequencing⁴ and combinations of above mentioned technologies⁵.

| NGS Dataset | NGS N50 (Mbp) | Total Scaffold Size (Mbp) | Scaffold N50 (Mbp) | # of NGS Anchored | Total NGS Anchored/Total NGS Length in Mbp (%) |
|---------------|---------------|---------------------------|--------------------|-------------------|--|
| Bionano Maize | 1.185 | 2,120 | 100 | 2,809 | 99.5% |
| Kakapo | 4.34 | 1,176 | 71 | 1,898 | 95.9% |
| Blackbird | 1.47 | 1,018 | 42 | 977 | 95.0% |

Table 2: Contiguity of sequence assemblies of 3 organisms is shown before (NGS N50) and after scaffolding (Scaffold N50) with Bionano maps built using DLS. Total assembly size and percentage of NGS sequence incorporated in the assembly is shown as well.

Assembly Conflicts and Resolution

The Bionano hybrid scaffold pipeline detects and resolves chimeric joins. Chimeric joins are typically formed when short reads, molecules, or paired-end inserts are unable to span across long DNA repeats. The errors appear as conflicting junctions in the alignment between the Bionano map and NGS assemblies.

When the hybrid scaffold pipeline detects a conflict, it analyzes the single-molecule data that underlies a Bionano map and assesses which assembly was incorrectly formed. If the Bionano map has long molecule support at the conflict junction, the sequence contig is automatically cut, removing the putative chimeric join (Figure 3). If it does not have strong molecule support, then the Bionano map is automatically cut. Both assemblies must have coverage spanning both sides of a chimeric join to detect and resolve these conflicts.

Hi-C-based scaffolding tools have become popular because of their ability to link even the shortest sequence contigs into chromosome-length scaffolds. However, this highly stochastic method based on the statistical interpretation of the frequency of chromatin cross-links within the nucleus shows a high number of errors in order and orientation of the contigs. Bionano maps can help identify and correct these errors^{11,15}, as shown in Figure 2.

Users can manually inspect all conflict resolution results. Bionano Solve notes the IDs and coordinates of the sequences and maps where conflicts have been detected and the corresponding resolution approaches taken. The scaffold can be edited in Bionano Access™ and modified, and then run again in the hybrid scaffold pipeline to produce a new set of scaffolds based on the manual conflict resolution. This manual enhancement process can be performed multiple times, giving users fine control in generating high-quality, complete hybrid scaffolds.

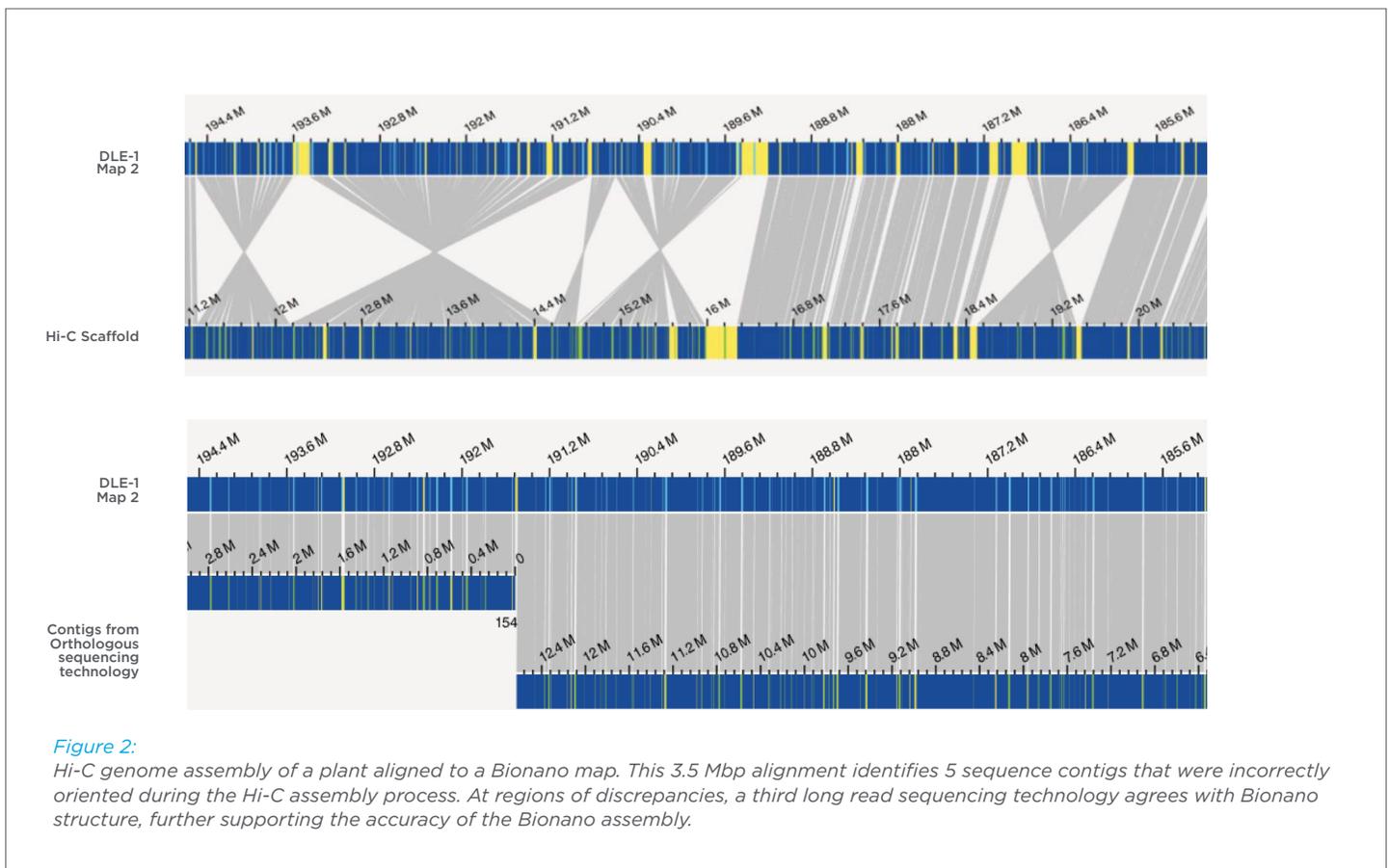


Figure 2: Hi-C genome assembly of a plant aligned to a Bionano map. This 3.5 Mbp alignment identifies 5 sequence contigs that were incorrectly oriented during the Hi-C assembly process. At regions of discrepancies, a third long read sequencing technology agrees with Bionano structure, further supporting the accuracy of the Bionano assembly.

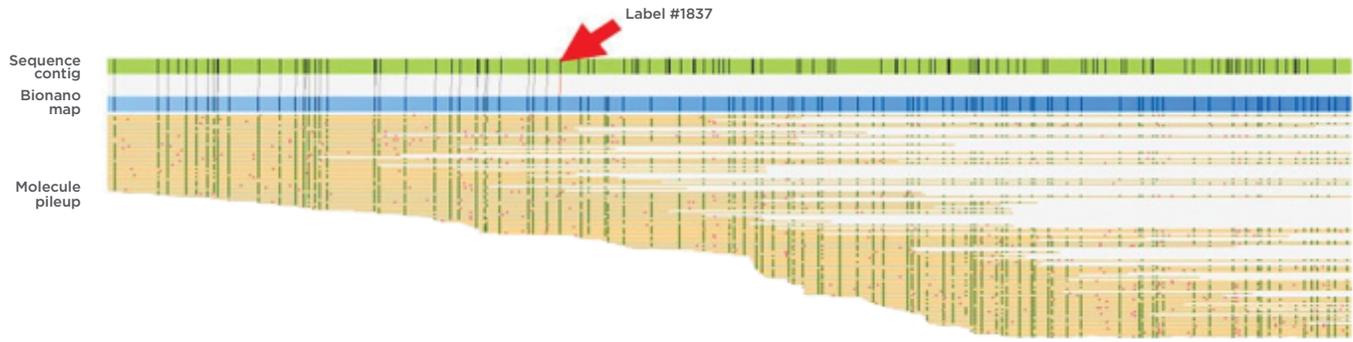


Figure 3:

Example of a conflict between a sequence contig and a Bionano map. (Top) The conflict junction is shown by the red arrow in the alignment between the sequence contig and the Bionano map. There is strong molecule support spanning the junction region on the genome map, so the sequence contig is cut at the label indicated.

Accuracy

A more accurate assembly doesn't just have better contiguity and fewer errors but is more functional as well. Genes and their regulatory sequences need to be assembled, ordered and oriented correctly to allow for a meaningful functional analysis. Figure 4 illustrates a region containing a muscle skeletal receptor kinase

(MuSK) gene in hummingbird—which may be of biological significance for its extreme flight skills. The PacBio *de novo* assembly failed at the dense repeats in the gene, splitting between two sequence contigs and failing to measure repeat array copy number. In contrast, Bionano hybrid scaffolding correctly brings the two pieces together to create one functional gene.

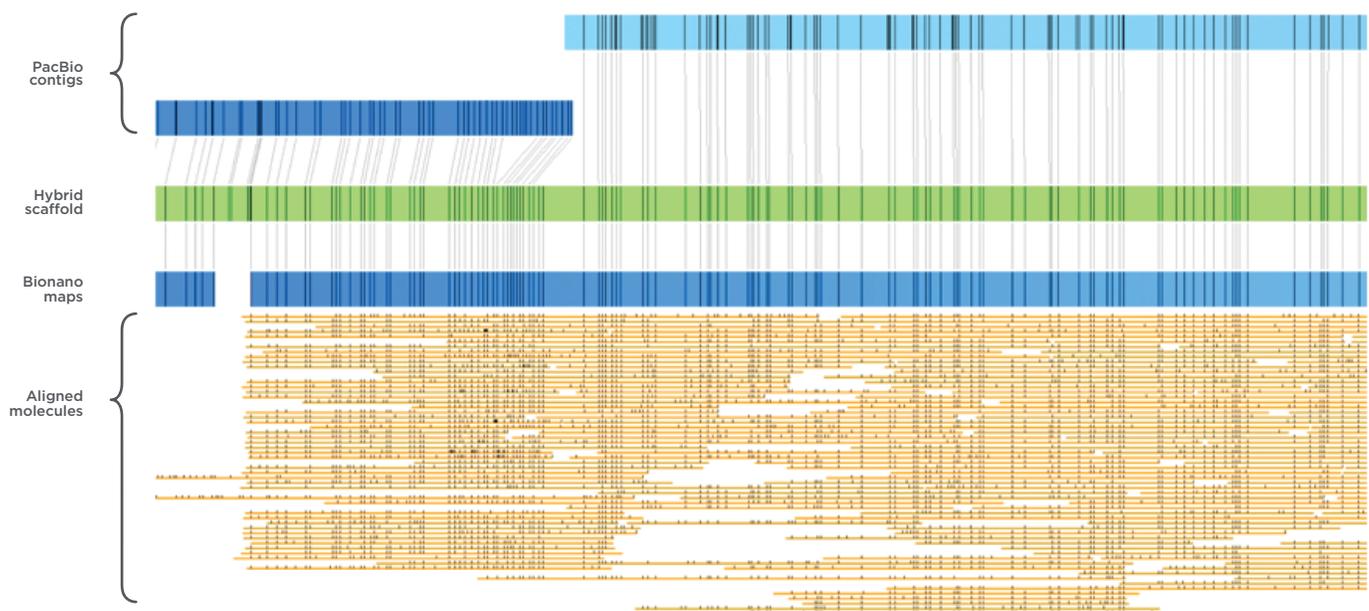


Figure 4:

The hummingbird Muscle Skeletal Receptor Kinase is split between two PacBio sequence contigs. Bionano hybrid scaffolding uses the genome map to anchor and orient the contigs, leading to a functional gene. Vertical lines represent *Nt.BspQI* recognition sites.

Higher Levels of Contiguity Using Two-Enzyme Hybrid Scaffolding

Assembly contiguity can be further increased by performing hybrid scaffolding with maps using two separate labeling enzymes. Two sets of Bionano maps, each generated with a different labeling enzyme, can be integrated with NGS sequences together. The two maps can both be generated using NLRS or could be one DLS map combined with one NLRS map. This integration enables the NGS sequences to function as a bridge to merge single-enzyme Bionano maps into two-enzyme maps that contain the sequence motif patterns from both labeling enzymes. Since the Bionano maps are generated independently they serve as orthogonal sources of evidence to detect and correct assembly errors in input data. The complementarity of different data also greatly improves the contiguity of the merged Bionano map while doubling the information density, which substantially increases the ability to anchor short NGS sequences in the final scaffolds.

The two-enzyme approach was validated on the human NA12878 genome, a model data set for which sequence data is publicly available. Three different assemblies were tested: two Discover assemblies of Illumina 250 bp pair-end sequence with different quality (Illumina-D2, N50: 0.08 Mbp, Illumina-D, N50: 0.18 Mbp); and PacBio, N50: 0.9 Mbp 46x with mean read length of 3.6 kbp. In each case, the assembly contiguity is higher with DLS alone than when combining two nicking endonuclease, but the two-enzyme approach combining DLS with NLRS maps improves the scaffold contiguity up to 1000-fold when compared to input NGS (Figure 5), anchors more sequence contigs in the final scaffolds and corrects more assembly errors in NGS sequences. The pipeline performs robustly in both animal and plant genomes and this approach greatly expands the type of NGS data that can be integrated with Bionano maps to produce highly accurate and contiguous assemblies for complex genomes.

The two-enzyme scaffolding method improves the error correction even further. Since the Bionano maps were generated independently they serve as orthogonal sources of evidence to detect and correct assembly errors in input data.

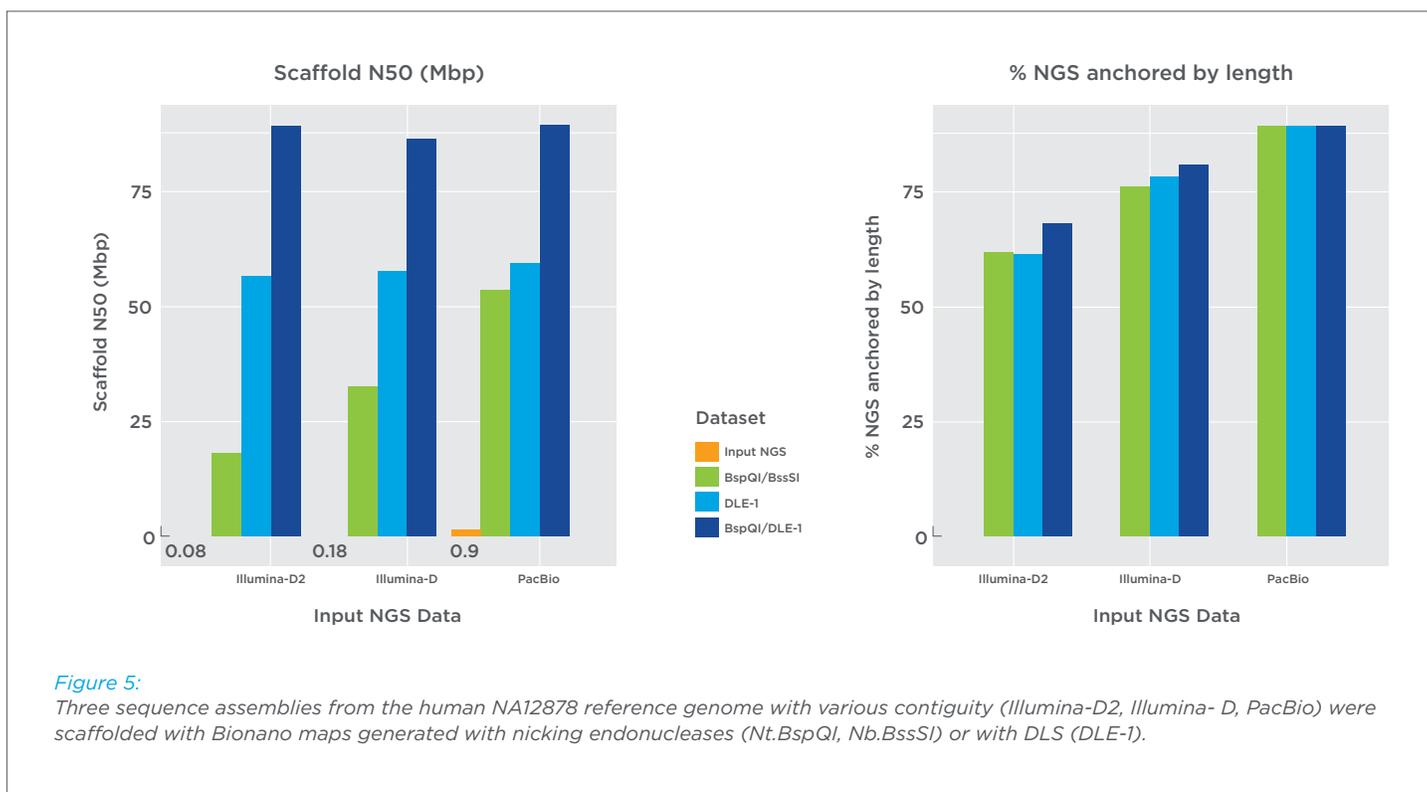


Figure 5: Three sequence assemblies from the human NA12878 reference genome with various contiguity (Illumina-D2, Illumina-D, PacBio) were scaffolded with Bionano maps generated with nicking endonucleases (*Nt.BspQI*, *Nb.BssSI*) or with DLS (*DLE-1*).

Third-party software

Several third-party software algorithms offer functionality beyond Bionano Access. A team from the Innovation Academy for Seed Design, Chinese Academy of Sciences published on their genome assembly method HERA (Highly Efficient Repeat Assembly), which resolves repeats efficiently by constructing a connection graph from an overlap graph¹². HERA was tested on rice, maize, human, and Tartary buckwheat genomes with single-molecule sequencing and Bionano optical mapping data. HERA correctly assembles most of the previously unassembled regions, resulting in dramatically improved, highly contiguous genome assemblies with newly assembled gene sequences.

A recent paper authored by scientists from the University of California, Riverside¹³ introduces a novel scaffolding algorithm called OMGS that can take advantages of multiple optical maps. OMGS solves several optimization problems to generate scaffolds with optimal contiguity and correctness. They use extensive experimental results to demonstrate that their tool outperforms existing methods when multiple optical maps are available.

The Gyde Bioinformatics tool¹⁴ contains hybrid assembly algorithms to simultaneously use sequencing, single-molecule and Bionano data to construct pseudomolecules representing entire chromosomes with cross-validation of all the supporting data. The tool includes a proprietary optical aligner and assembler, tools for assembling sequencing data, and a visualization/editing solution.

Phasing of parental alleles during assembly

Generating complete diploid assemblies with correct phasing has always been a challenge. Trio-binning of proband sequence reads using parental reads has been applied recently to partition proband reads prior to the assembly of the binned paternal and maternal sets¹⁶. This approach has been used in Bionano optical mapping data as well, presented as a poster at PAG 2020¹⁵. Trio-binning of proband Bionano molecules to the Bionano parental haplotype-aware assemblies gives contiguous assembly of each allele in a human trio dataset, and several animal genomes.

Cost Considerations

Bionano hybrid scaffolding makes an assembly better for a low cost. Adding a Bionano genome map to your assembly costs as little as \$500 in materials for up to human size genomes and remains affordable for larger genomes as well (e.g. wheat). These costs are significantly lower than investments required for PacBio sequencing, Dovetail or NRGene assemblies. No matter what your sequencing strategy is, adding Bionano to your assembly will help achieve significant improvements in contiguity and accuracy, produce a better assembly and thus a superior publication at a reasonable cost.

Alternatively, the improvements in contiguity using Bionano hybrid scaffolding allow you to reduce the sequencing coverage necessary to produce an assembly of a certain quality. A low coverage 30x PacBio assembly scaffolded with Bionano maps typically produces an assembly of superior contiguity comparable to 80x coverage with PacBio alone. Depending on the organism's genome size, a significantly lower coverage PacBio sequencing coupled with Bionano maps can reduce the cost by tens of thousands of dollars – far more than the cost to generate Bionano data.

Discussion

Combining NGS and Bionano optical mapping data produces assemblies of the highest quality. This approach offers an affordable solution to improve fragmented draft assemblies and build the highest-quality assemblies containing accurate long-range information.

Bionano hybrid scaffolding is agnostic to the sequence technology used. Recent publications have scaffolded assemblies based on Illumina sequencing alone, PacBio alone, 10x Genomics assemblies, Oxford Nanopore Assemblies, NRGene assemblies, nanopore sequencing, and combinations of those.

| | AK1 | HX1 | NA12878 | NA12878 | NA24385 | GRCh38 |
|--|---------|---------|-------------------------|---------|--------------------|----------|
| Sequencing | PacBio | PacBio | Illumina + 10x Genomics | PacBio | PacBio | Sanger |
| Scaffolding | Bionano | Bionano | Bionano | Bionano | Bionano two-enzyme | multiple |
| Input N50 (Mbp) | 17.92 | 8.325 | 7.03 | 1.56 | 4.7 | 56.41 |
| Hybrid Scaffold N50 (Mbp) scaffold | 44.85 | 21.979 | 33.5 | 26.83 | 80.46 | 67.79 |
| Fold Improvement after Bionano hybrid scaffold | 2.5x | 2.6x | 4.8x | 17.2x | 17.1x | |

Table 3:
Assembly statistics for a number of human reference genomes (data from Genome in a Bottle Consortium^{2,3,9,10}).

Bionano maps can error correct input sequence assemblies. Any of the scaffolding technologies using synthetic long reads or DNA cross-linking provide some sort of error correction compared to short-read assemblies alone. However, since they are NGS based, they suffer from most of the same problems plaguing short-read only assemblies. Only Bionano optical Mapping provides non-sequencing based, orthogonal genome structure data in a high throughput way, allowing for a completely independent error correction.

Recent publications on reference genomes for wheat, banana, bed bug and maize^{5,6,7,8} all included Bionano data to create higher contiguity and/or correct assembly errors. All major human reference

genome publications use Bionano optical mapping data as well, including the NA12878 genome², the Chinese reference genome⁹ and the Korean reference genome¹⁰. The contiguity of these recent genome publications combining *de novo* sequence assemblies with Bionano maps approach or surpass that of the hGCR38 reference genome (Table 3). The Vertebrate Genome Project, The Darwin Tree of Life Project and many other large genome consortia use Bionano optical mapping data to produce the highest quality assemblies. **Including Bionano optical mapping data into *de novo* genome assemblies has become a gold standard.**



References: 1.J. W. Clouse et al The Amaranth Genome: Genome, Transcriptome, and Physical Map Assembly The Plant Genome (2016) 2.Pendleton, M., Sebra, R., et al. Assembly and diploid architecture of an individual human genome via single-molecule technologies. Nature Methods (2015); e3454 3.Mostovoy J et al. A hybrid approach for *de novo* human genome sequence assembly and phasing Nature Methods (2016) 4.Deschamps S et al. A chromosome-scale assembly of the sorghum genome using nanopore sequencing and optical mapping. Biorxiv 2018 5.Zimin et al Hybrid assembly of the large and highly repetitive genome of *Aegilops tauschii*, a progenitor of bread wheat, with the mega-reads algorithm bioRxiv 2016 6.Martin et al. Improvement of the banana “*Musa acuminata*” reference sequence using NGS data and semi-automated bioinformatics methods BMC Genomics (2016) 7.Rosenfeld et al. Genome assembly and geospatial phylogenomics of the bed bug *Cimex lectularius* Nature Communications (2016) 8.Dong et al Analysis of tandem gene copies in maize chromosomal regions reconstructed from long sequence reads PNAS (2016) 9.Shi et al Long-read sequencing and *de novo* assembly of a Chinese genome Nature Communications (2016) 10.Seo JS et al *de novo* assembly and phasing of a Korean human genome Nature (2016) 11. Udall JA, Dawe RK. Validating Genome Assemblies by Optical Mapping. Plant Cell. (2018) 12. Du, H., Liang, C. Assembly of chromosome-scale contigs by efficiently resolving repetitive sequences with long reads. Nat Commun 10, 5360 (2019) 13. Pan W., Jiang T., Lonardi S. OMGs: Optical Map-Based Genome Scaffolding. In: Cowen L. (eds) Research in Computational Molecular Biology, RECOMB 2019. Lecture Notes in Computer Science, vol 11467. Springer, Cham (2019) 14. Rigault, Philippe & Dimech, Adam. Optical and physical mapping with local finishing enables megabase-scale resolution of agronomically important regions in the wheat genome. Genome Biology. (2018) 15. Lee J Identifying Parental Alleles with Bionano Genomics’ Ultra-High Molecular Weight DNA. Poster, PAG (2020) 16. Koren S, Rhie A, Walenz BP, Dilthey AT, Bickhart DM, Kingan SB, Hiendleder S, Williams JL, Smith TPL, Phillippy AM. *De novo* assembly of haplotype-resolved genomes with trio binning. Nat Biotechnol. (2018)

For general information about the Saphyr™ System, please contact info@bionanogenomics.com or visit bionanogenomics.com

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