

# Hybrid Scaffolding Improves Genome Assembly Accuracy and Contiguity

## Next-Generation Mapping Reveals True Long Range Structure of the Genome and Reduces Sequencing Costs

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

**The genomes of most higher organisms are highly repetitive.** Two thirds of human and most mammalian genomes consist of repeats, and many plant genomes have even higher repeat content. Short reads usually fail to span long 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 and limited long-range information. NGS technology is based on fragmenting DNA molecules, reading just hundred(s) of basepairs, and using algorithms to reassemble these fragments.

**The introduction of long-read sequencing has led to improved assembly contiguity and accuracy** as well, but can be time-consuming and expensive, especially when deep coverage or the spanning of long tandem repeats is required. Read lengths are still limited to tens of kilobasepairs.

**Recently, synthetic long-read technologies like that of 10x Genomics have gained momentum.**

Using a barcoding method to link short reads, some mid-range structural information is retained, thus improving the contiguity of NGS assemblies. Synthetic long reads are still plagued by some challenges inherent to NGS technology. These challenges include failure to disambiguate interspersed repeat units and correctly assemble and size long repeat array; assembly gaps due incomplete coverage and GC bias in PCR amplification and sequencing; and lack of long range structural information caused by fragmenting of the DNA and reads that are too short to span and correctly resolve larger structural variation.

**Only extremely long molecules, ranging in size from hundreds of thousands to millions of base pairs, provide accurate structure of the genome.** Bionano Genomics' Next-Generation Mapping (NGM) visualizes long DNA molecules in their native state. Long range genomic structural information is preserved and directly 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 order and orientation of sequence elements in the map can be used for anchoring NGS contigs and detecting structural variation.

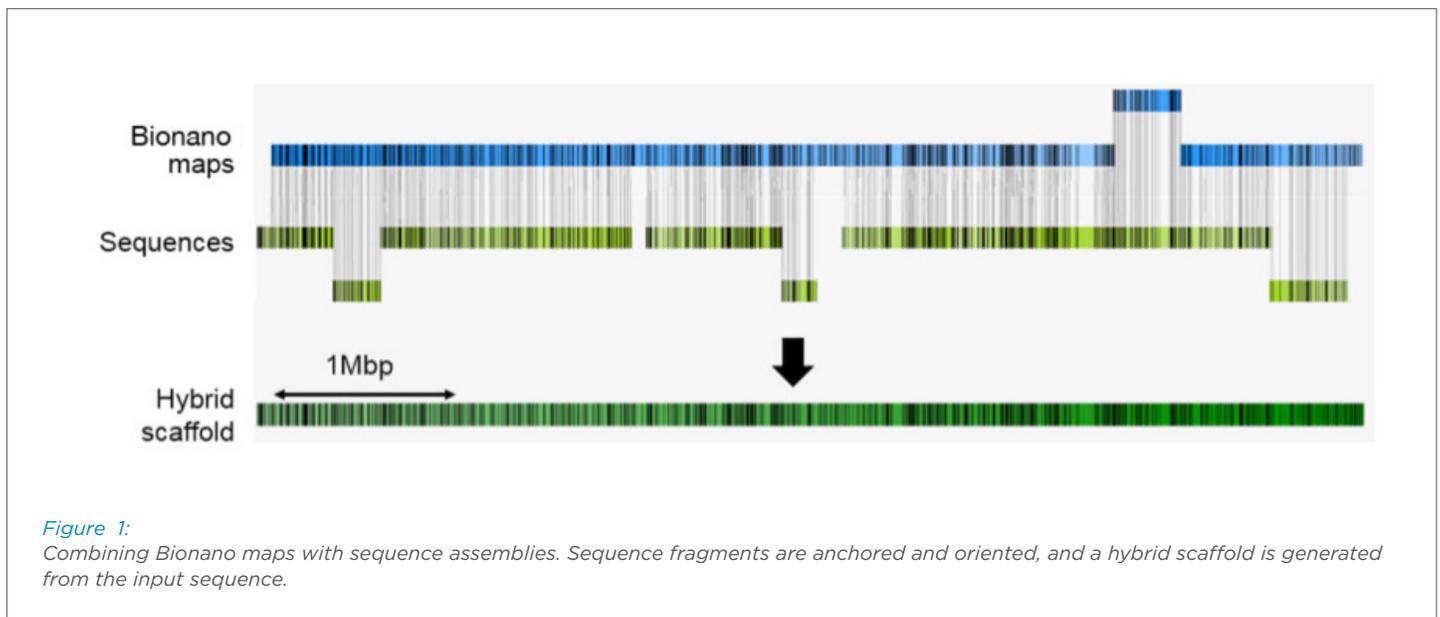
Megabase size molecules of genomic DNA are labeled, linearized and uniformly stretched in high density NanoChannel arrays, and imaged on the Bionano instrument. Using a nicking endonuclease, a specific 6 or 7 basepair sequence is labeled approximately 10 times per 100 kbp. The label patterns allow each long molecule to be uniquely identified and aligned. Using pairwise alignment of the single molecules, consensus genome maps are constructed, refined, extended and merged. Genome maps can be created using different endonucleases to generate broader coverage and higher label density.

## Hybrid Scaffold construction

The *de novo* Bionano genome 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 putative nick sites in the sequence based on the nicking endonuclease-specific recognition site. These *in silico* maps for the

sequence contigs are then aligned to the *de novo* Bionano genome 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 Hybrid Scaffold and viewing scaffolding results.



## Contiguity and Completeness

**The hybrid scaffolding process considerably reduces the number of contigs found in the initial NGS assembly**, improving assembly accuracy and quality while reducing the need for deep sequencing coverage.

The hybrid scaffolding approach can yield significant improvements in contiguity, as expressed by the assembly N50 values, across genomes as seen in Table 1. We created hybrid scaffolds for three genomes (human, goat, and maize). This process

improved contiguity by as much as 13-fold. The Bionano Solve pipeline makes near-complete use of the available input assemblies, taking into account more than 84% of the total length of the sequence and at least 93% of the genome map (Table 2).

**Bionano hybrid scaffolding is agnostic to the sequence technology used.** Recent publications featured scaffolded assemblies using Illumina sequencing alone<sup>1</sup>, PacBio alone<sup>2</sup>, 10x Genomics assemblies<sup>3</sup>, and combinations of those<sup>4</sup>.

*Table 1:  
Contiguity and genome coverage of hybrid scaffolds.*

Species	Sequence N50 (Mbp)	Bionano Map N50 (Mbp)	Hybrid Scaffold Contiguity		Hybrid Scaffold Coverage	
			Hybrid Scaffold N50 (Mbp)	N50 Fold Increase	Hybrid Scaffold Size (Mbp)	% of Known Reference Assembly
Human NA1287	0.90	3.92	11.94	13.33	2835.18	91.8
Goat	4.68	1.59	23.85	5.10	2643.93	104.7
Maize	1.04	2.47	10.00	9.63	2119.64	102.6

*Table 2:  
Use of input assemblies in scaffolding.*

Species	Amount of Sequence Data Utilized in Hybrid Scaffold (Mbp)	Amount of Bionano Data Utilized in Hybrid Scaffold (Mbp)
Human NA12878	2,576 (84.03%)	2,804.64 (98.07%)
Goat	2,498 (95.13%)	2,572.30 (93.60%)
Maize	2,008 (95.42%)	2,079.08 (98.11%)

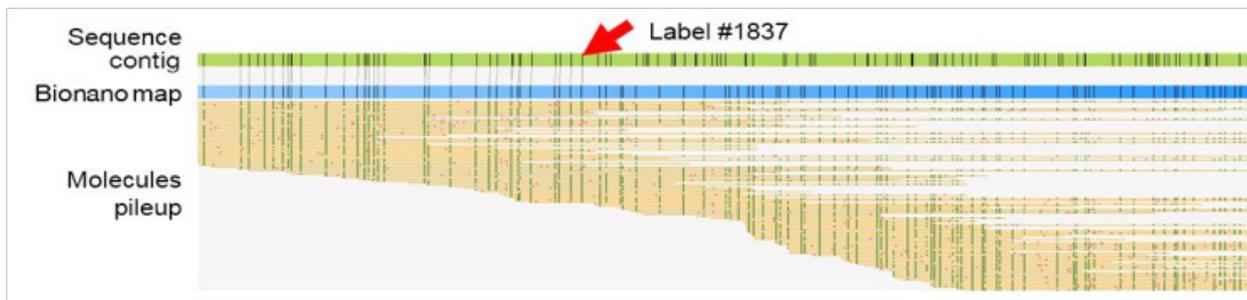
## Assembly Conflicts and Resolution

**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 2). 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.

Automated cuts using Bionano Solve help to resolve conflicts with a high level of accuracy. The majority of cuts made using Bionano Solve can be confirmed by comparison to the species' reference assembly (Table 3). There are several reasons why some cuts cannot be confirmed: the reference assembly is incomplete, the two separate input assemblies may represent different alleles, or the chimeric joins may have been caused by segmental duplications that are too long for Bionano molecules to resolve.

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. This file can be edited 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:** Example of a conflict between a sequence contig and a Bionano map. The conflict junction as 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 is cut at the label indicated.

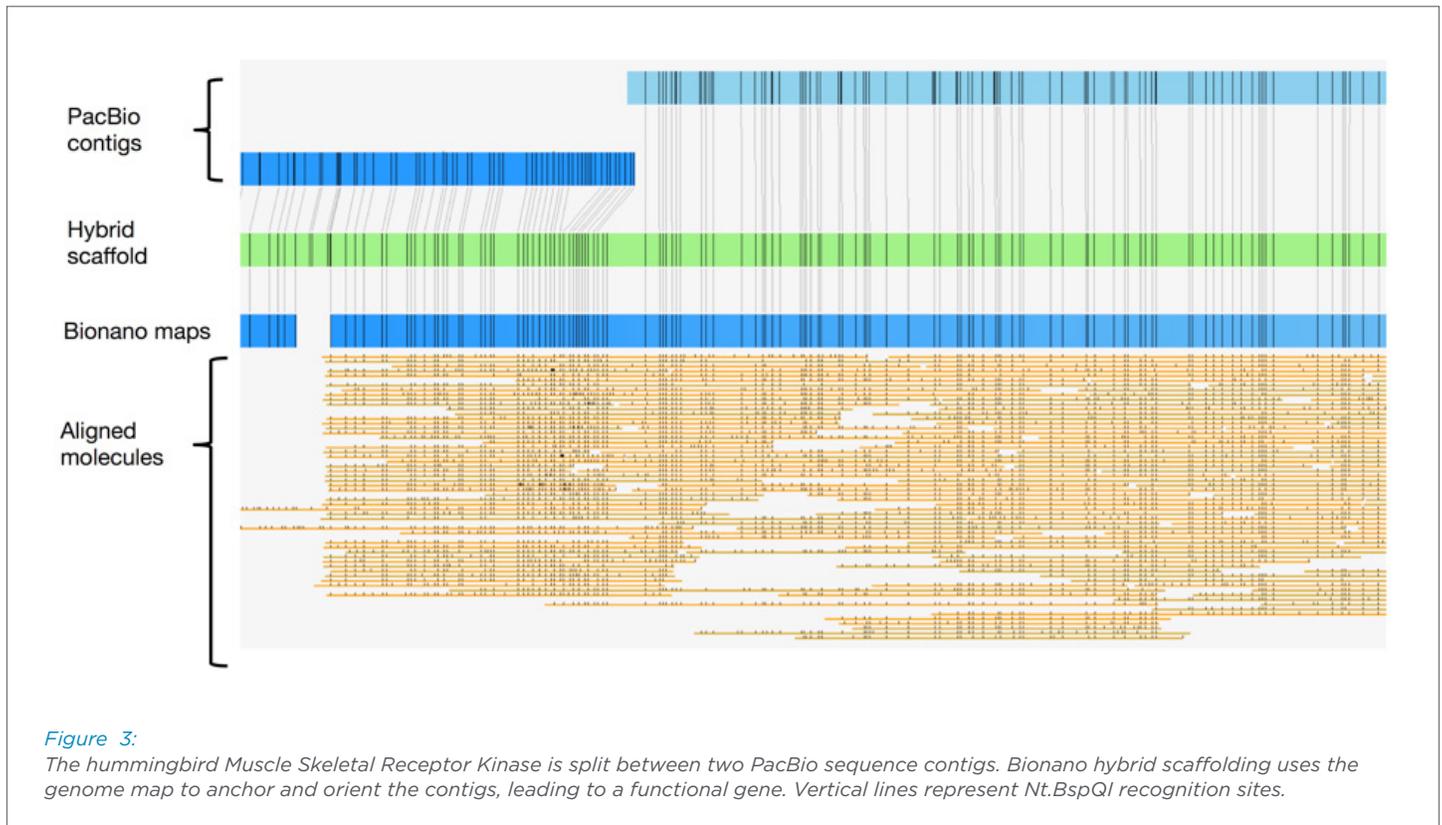
**Table 3:** The number of cuts performed by the hybrid scaffold pipeline to the assemblies and the number of cuts confirmed when aligned to references.

Species	# of Cuts on Sequence Confirmed/Total	# of Cuts on Bionano Confirmed/Total
Human NA12878	4 / 6 (67%)	1 / 1 (100%)
Goat	66 / 79 (84%)	11 / 16 (69%)
Maize	24 / 26 (92%)	12 / 13 (92%)

## 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 3 illustrates a region containing a muscle skeletal receptor kinase (MuSK) gene in hummingbird – which may be of

biological significance for the extreme flight skills of this bird. The *de novo* PacBio assembly failed at the dense repeats in the gene, leading to it being split between two sequence contigs and failing to measure repeat array copy number. Bionano hybrid scaffolding correctly brings the two pieces together to create one functional gene.



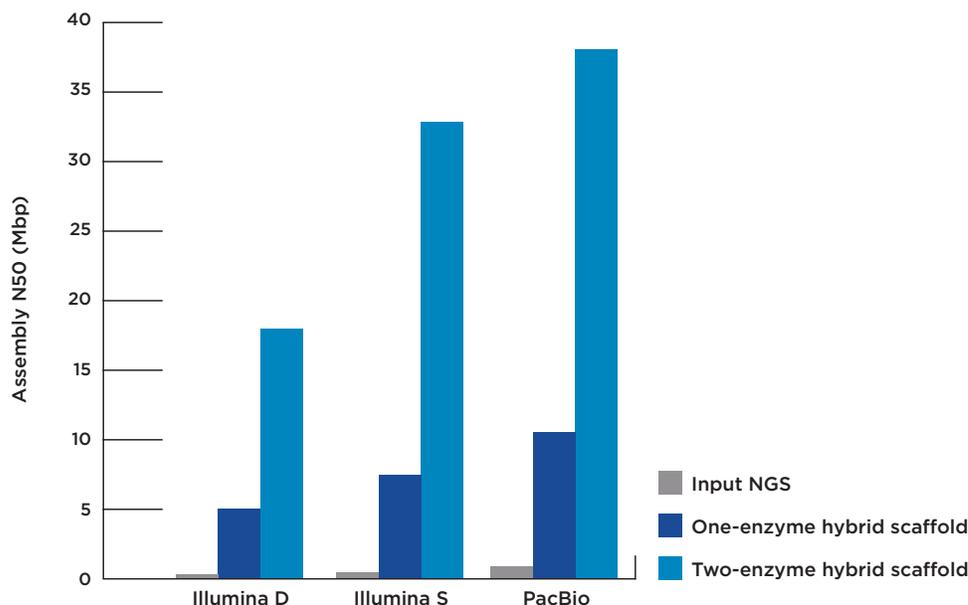
## Higher Levels of Contiguity Using Two-Enzyme Hybrid Scaffolding

**Assembly contiguity can be further increased by performing hybrid scaffolding with maps using two separate nicking enzymes.** Two sets of Bionano maps, each generated with a different nicking enzyme, can be integrated with NGS sequences together. This 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 nicking 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: Illumina-D, 51x of 250 bp pair-

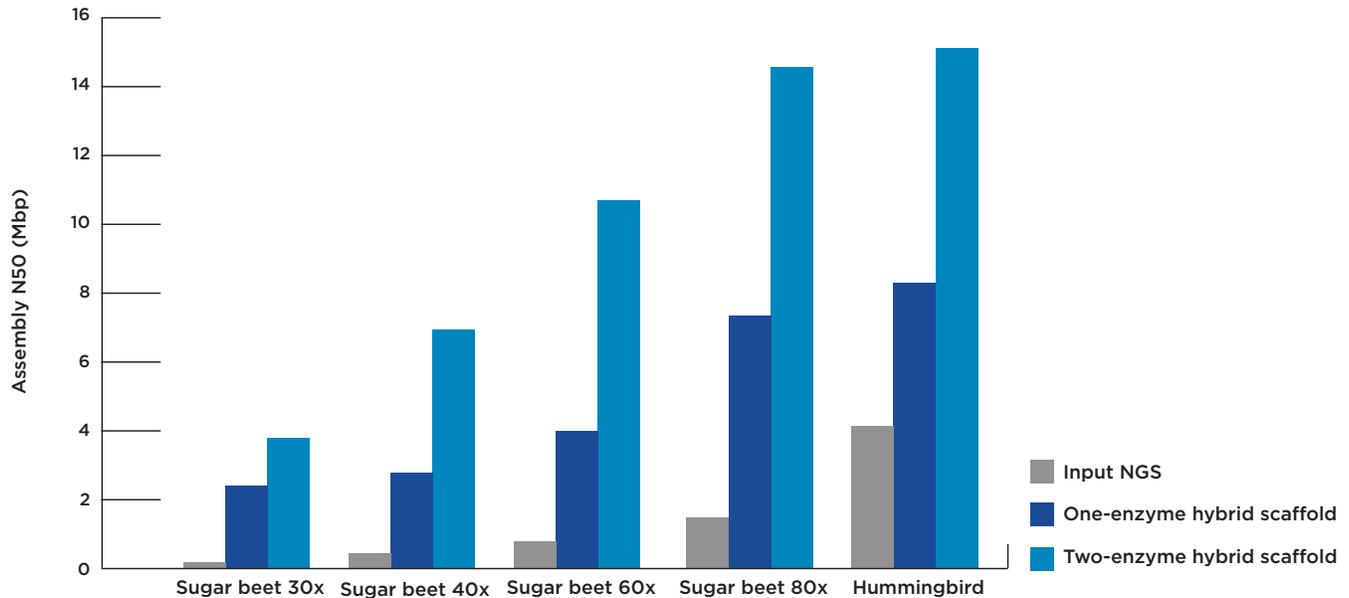
end sequence; Illumina-S, 40x of 101 bp pair-end and 25x of 2.5-2.5 kbp mate-pair sequence; and PacBio, 46x with mean read length of 3.6kbp. Compared to the published single-enzyme hybrid-scaffolds, the two-enzyme approach improves the scaffold contiguity 3-fold (up to 100-fold when compared to input NGS, Figure 4), anchors 30% more sequence contigs in the final scaffolds and corrects 50% more assembly errors in NGS sequences. The pipeline performs robustly in both animal and plant genomes as well (Figure 5). 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 evidences to detect and correct assembly errors in input data. Compared to the published one-enzyme hybrid-scaffolds, the two-enzyme approach corrects up to 50% more assembly errors in NGS sequences.



**Figure 4:**

Improvements in NA12878 assembly contiguity after hybrid scaffold with one-enzyme and two-enzyme genome maps. Illumina-D: 51x of 250 bp pair-end sequence; Illumina-S: 40x of 101 bp pair-end and 25x of 2.5-2.5 kbp mate-pair sequence; PacBio: 46x with mean read length of 3.6kbp



**Figure 5:**

Improvements in sugar beet and hummingbird assembly contiguity after hybrid scaffolding with Bionano genome maps using one-enzyme and two-enzymes. For sugar beet, the fold coverage of the PacBio de novo assemblies is shown.

## Cost considerations

**Bionano hybrid scaffolding makes an assembly better for a low cost.** Adding a Bionano genome map to your assembly costs less than \$1000 in materials for smaller genomes, and remains affordable for larger genomes as well. This compares extremely favorably to PacBio sequencing, Dovetail or NRGene assemblies. No matter what your sequencing strategy is, adding Bionano to your assembly is a good decision. The 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. As the example of sugar beet shows in Figure 5, a 30x PacBio assembly scaffolded with Bionano maps produces an assembly of superior contiguity than 80x PacBio alone. Depending on the organism's genome size, a

significant reduction in PacBio sequencing can reduce the cost by tens of thousands of dollars – far more than the cost to generate the Bionano data.

## Discussion

### Combining NGS and Bionano NGM 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, NRGene assemblies, and combinations of those.

**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 NGM 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<sup>4,5,6,7</sup> all included Bionano data** to create higher contiguity and/or correct

assembly errors. All major human reference genome publications used NGM data as well, including the NA12878 genome<sup>2</sup>, the Chinese reference genome<sup>8</sup> and the Korean reference genome<sup>9</sup>. 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 4). **Including Bionano mapping data into *de novo* genome assemblies has become a *de facto* standard.**

*Table 4:  
Assembly statistics for a number of human reference genomes<sup>2,3,8,9</sup>*

	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	

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.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 5.Martin et al. Improvement of the banana "Musa acuminata" reference sequence using NGS data and semi-automated bioinformatics methods BMC Genomics (2016) 6.Rosenfeld et al. Genome assembly and geospatial phylogenomics of the bed bug Cimex lectularius Nature Communications (2016) 7.Dong et al Analysis of tandem gene copies in maize chromosomal regions reconstructed from long sequence reads PNAS (2016) 8.Shi et al Long-read sequencing and de novo assembly of a Chinese genome Nature Communications (2016) 9.Seo JS et al de novo assembly and phasing of a Korean human genome Nature 2016

For general information about the Irys® System, please contact [info@bionanogenomics.com](mailto:info@bionanogenomics.com) or visit [bionanogenomics.com](http://bionanogenomics.com)