

# BIOINFORMATICS APPLICATION NOTE #1

## Benchmark Results

### OmniTier CompStor Novos™ Alignment shows greater variant calling accuracy than GATK Best Practices across the seven GIAB datasets with accelerated runtimes

OmniTier Inc. | 1591 McCarthy St, Milpitas, CA 95035 | 2720 Superior Dr NW, Rochester, MN, 55901

#### SUMMARY

OmniTier's CompStor Novos™ bioinformatics appliance shows greater variant calling accuracy across the seven Genome in a Bottle (GIAB) datasets than GATK Best Practices pipeline as measured by F<sub>1</sub> scores, minimum error counts, and ROC curves.

The CompStor Novos™ appliance demonstrates accelerated run-times, completing variant calling in 2 hours on a 2-node configuration versus 9.4 hours for the GATK Best Practices pipeline. Further runtime reduction below 1 hour is feasible with more nodes.

## 1 INTRODUCTION

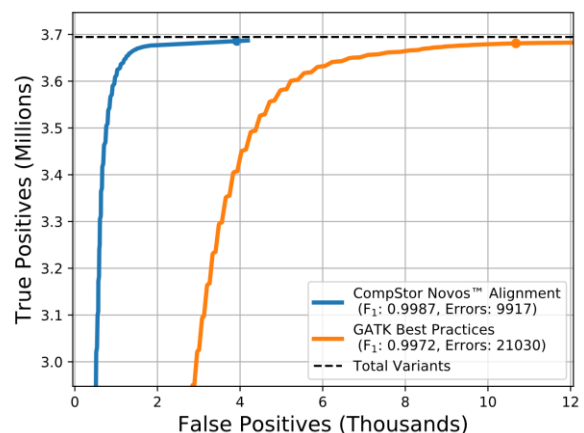
Accurate and fast variant calling in whole genome sequencing (WGS) drives precision medicine adoption. CompStor Novos™ is a scalable compute appliance utilizing a tiered-memory architecture of DRAM and SSDs to accelerate and improve variant calling in germline WGS pipelines. It is a *dual pipeline* utilizing *de novo* assembly and alignment techniques. This strategy increases the opportunity to find unique variants while driving down runtimes. The appliance is controlled by an intuitive web application interface and enables batch processing automation.

For variant calling, OmniTier uses a domain-optimized deep learning methodology to produce fewer false positives and more true positives in germline WGS. This improved performance is demonstrated across the 7 GIAB datasets made available by the National Institute of Standards and Technology (NIST)<sup>1</sup>.

## 2 COMPARING DATA

GIAB datasets (HG001-HG007) were downloaded from NIST GIAB website<sup>2</sup>. The GATK 4.0.1.2 pipeline software was downloaded from the Broad Institute website<sup>3</sup>. F<sub>1</sub> scores are calculated from recall (fraction of true variants detected) and precision (fraction of variants called that are true):

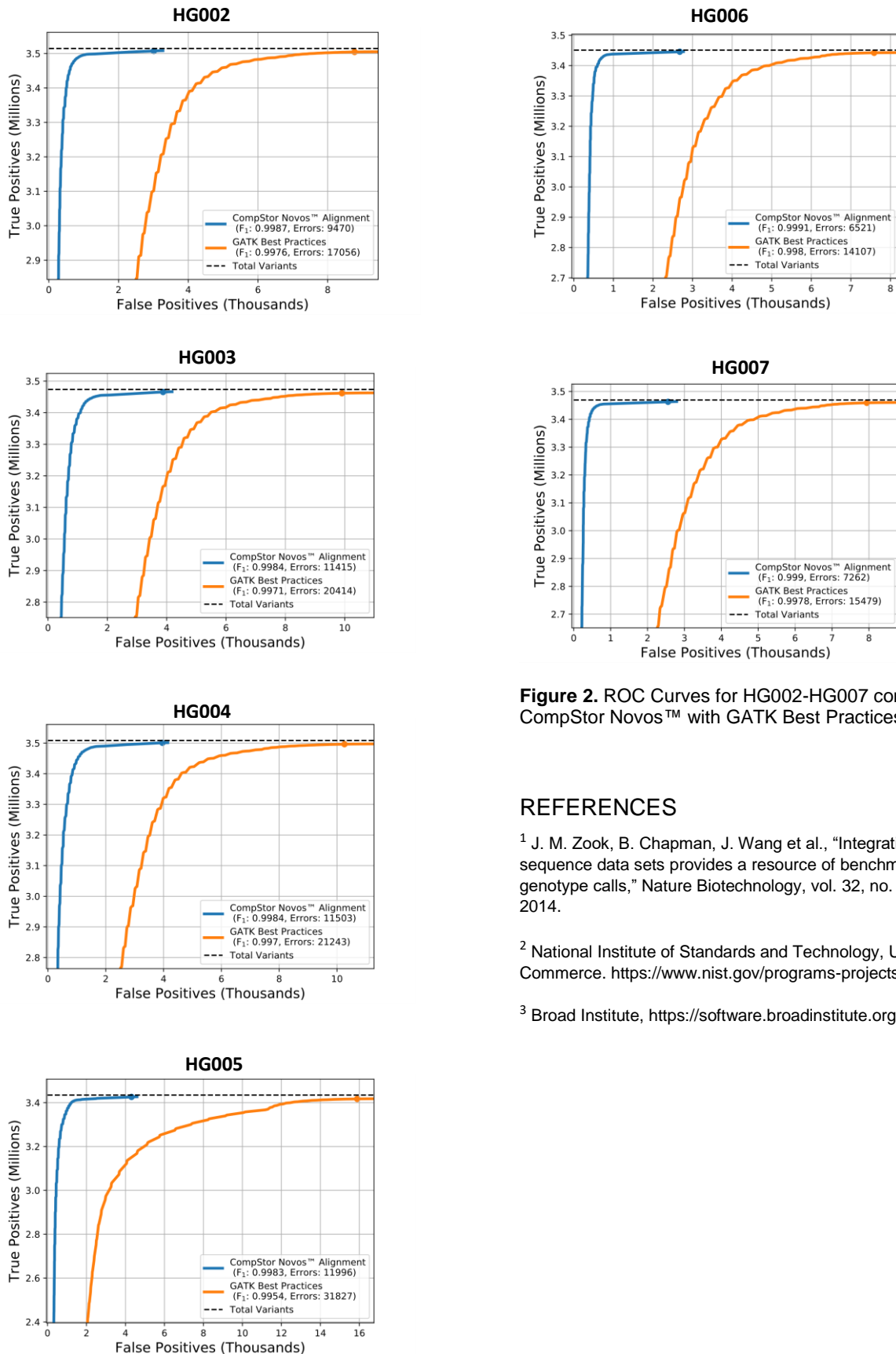
$$F_1 = \frac{2}{(\text{recall}^{-1} + \text{precision}^{-1})}$$



**Figure 1.** Receiver Operator Characteristics (ROC) Curve comparing HG001 true positives and false positives between CompStor Novos™ Alignment and GATK Best Practices

GIAB Dataset	CompStor Novos™	GATK Best Practices	Percentage Improvement
HG001	9,917	21,030	52.8%
HG002	9,470	17,056	44.5%
HG003	11,415	20,414	44.1%
HG004	11,503	21,243	45.9%
HG005	11,996	31,827	62.3%
HG006	6,521	14,107	53.8%
HG007	7,262	15,479	53.1%
<b>Total Errors Avg</b>	<b>9,726</b>	<b>20,165</b>	<b>51.8%</b>
<b>F<sub>1</sub> Average</b>	<b>.9987</b>	<b>.9972</b>	

**Table 1.** Total Errors and F<sub>1</sub> Scores comparing CompStor Novos™ with GATK Best Practices for all NIST GIAB datasets



**Figure 2.** ROC Curves for HG002-HG007 comparing CompStor Novos™ with GATK Best Practices

## REFERENCES

- 1 J. M. Zook, B. Chapman, J. Wang et al., "Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls," Nature Biotechnology, vol. 32, no. 3, pp. 246–251, 2014.
- 2 National Institute of Standards and Technology, US Department of Commerce. <https://www.nist.gov/programs-projects/genome-bottle>.
- 3 Broad Institute, <https://software.broadinstitute.org/gatk/>.